Bioinformatics



Multiple Alignment, Patterns & Profiles

David Gilbert

Bioinformatics Research Centre

www.brc.dcs.gla.ac.uk

Department of Computing Science, University of Glasgow

Lecture summary

- Characterising *families* of sequences
- Multiple sequence alignment
- Weight matrices
- Searching for distant relatives: beyond Blast PSI-Blast
- Patterns
- Pattern discovery
- Rating & using patterns

Multiple Sequence Alignment

- Why do MSA?
 - Help prediction of the secondary and tertiary structures of proteins of new sequences
 - Help to find motifs or signatures characteristic of protein family

VTISCTGSSSNIGAG-NHVKWYQQLPG VTISCTGTSSNIGS--ITVNWYQQLPG LRLSCSSSGFIFSS--YAMYWVRQAPG LSLTCTVSGTSFDD--YYSTWVRQPPG PEVTCVVVDVSHEDPQVKFNWYVDG--ATLVCLISDFYPGA--VTVAWKADS--AALGCLVKDYFPEP--VTVSWNSG---

MSA

VTISCTGSSSNIGAG-NHVKWYQQLPG VTISCTGTSSNIGS--ITVNWYQQLPG LRLSCSSSGFIFSS--YAMYWVRQAPG LSLTCTVSGTSFDD--YYSTWVRQPPG PEVTCVVVDVSHEDPQVKFNWYVDG--ATLVCLISDFYPGA--VTVAWKADS--AALGCLVKDYFPEP--VTVSWNSG---VSLTCLVKGFYPSD--IAVEWWSNG--

- 8 fragments from immunoglobulin sequences
- alignment highlights
 - conserved residues,
 - -conserved regions

-more sophisticated patterns, like the dominance of hydrophobic residues (V,L,I) at fragment positions 1 and 3.

- http://www.techfak.uni-bielefeld.de/bcd/Curric/MulAli

MSA

VTISCTGSSSNIGAG-NHVKWYQQLPG VTISCTGTSSNIGS--ITVNWYQQLPG LRLSCSSSGFIFSS--YAMYWVRQAPG LSLTCTVSGTSFDD--YYSTWVRQPPG PEVTCVVVDVSHEDPQVKFNWYVDG--ATLVCLISDFYPGA--VTVAWKADS--AALGCLVKDYFPEP--VTVSWNSG---VSLTCLVKGFYPSD--IAVEWWSNG--

•The alignment can also enable us to infer the evolutionary history of the sequences.

• It looks like the first 4 sequences and the last 4 sequences are derived from 2 different common ancestors, that in turn derived from a "root" ancestor.

- But true phylogentic analysis is more complex
- http://www.techfak.uni-bielefeld.de/bcd/Curric/MulAli

Multiple sequence aligment - methods

- Simultaneous: N-wise alignment (adapted from *pairwise* approach)
 - uses N-dimension dynamic programming matrix.
 - Complexity is for global alignment
 - $O(m_1m_2)$ [2 sequences length $m_1 \& m_2$]
 - $O(m^2)$ [2 sequences of length m]
 - $O(m^n)$ [n sequences of length m]
 - Ten sequences of length 1000 requires $1000^{10} = 10^{\circ}$
 - Approximate age of universe in pico-seconds
 - Combinatrial explosion!
 - Thus only good for short sequences.
- Manual (!)
- Heuristic...

Multiple sequence aligment - methods

- Heuristic methods, e.g. Progessive -- ClustalW:
 - Split multiple alignment into pairwise alignments (?how?)
 - optimise locally greedy at each step
- Many possibilities as to how the sequence of (pairwise) alignments can be built
- Must attempt to minimise errors introduced in early alignments which will accumulate during the progressive alignment
- Can be achieved in part by aligning the MOST similar sequences in turn
- Employ a phylogenetic tree to 'guide' the progressive alignment
 - compute pairwise sequence identities
 - construct binary tree (can output phylogenetic tree)
 - align similar sequences in pairs, add distantly related ones later.
- No guarantee that the global optimum will be found
 - But provides a computationally tractable and biologically useful algorithm

Multiple Sequence Alignment

- Outline of CLUSTAL (Thomson et al 1994)
 - Calculate the pairwise similarity scores for the sequences
 - Can use full dynamic programming approach
 - Employing similarity score create a phylo tree (UPGMA)
 - From tree produce weights for each sequence
 - Based on similarities
 - High weighting to dissimilar sequences
 - Low weighting to similar sequences
 - Weighting used when combining alignments
 - Employing tree structure as a guide perform progressive pairwise alignments

Multiple Sequence Alignment



Multiple sequence alignment (globins)

CLUSTAL W (1.81) multiple sequence alignment

Human VHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKV 60

Gorilla VHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKV 60

Rabbit VHLSSEEKSAVTALWGKVNVEEVGGEALGRLLVVYPWTQRFFESFGDLSSANAVMNNPKV 60

Human KAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGK 120

Gorilla KAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFKLLGNVLVCVLAHHFGK 120

- Rabbit KAHGKKVLAAFSEGLSHLDNLKGTFAKLSELHCDKLHVDPENFRLLGNVLVIVLSHHFGK 120
- Pig KAHGKKVLQSFSDGLKHLDNLKGTFAKLSELHCDQLHVDPENFRLLGNVIVVVLARRLGH 120
- Human EFTPPVQAAYQKVVAGVANALAHKYH 146
- Gorilla EFTPPVQAAYOKVVAGVANALAHKYH 146
- Rabbit EFTPOVOAAYOKVVAGVANALAHKYH 146
- Pig DFNPNVQAAFQKVVAGVANALAHKYH 146
 - * * ****• *************



Multiple alignments

• Analyse gene families

- reveal (subtle) conserved family characteristics



(c) David Gilbert 2007

Profile (frequency matrix)

У	d	G	G			V	e	A	⊥
	_	~	~		* 7 *	T 7	-	7	٦
Y	Ε	G	G	A	V	V	Q	A	L
F	D	-	G	I	L	V	Q	А	V
F	Ε	G	G	I	L	V	Ε	А	L
Y	D	G	G	-	_	_	E	А	L
Y	D	G	G	A	V	_	E	A	L
1	2	3	chara 4	5	6	7	8	9	10
	1 Y Y F F Y	1 2 Y D Y D F E F D Y E	1 2 3 Y D G Y D G F E G F D - Y E G	$\begin{array}{c} chara \\ 1 \\ 2 \\ 3 \\ 4 \\ \end{array}$ $\begin{array}{c} Y \\ D \\ G \\ G \\ F \\ E \\ G \\ G \\ \end{array}$ $\begin{array}{c} G \\ G \\ G \\ G \\ \end{array}$	$\begin{array}{c} \text{characters} \\ 1 & 2 & 3 & 4 & 5 \\ \end{array}$ $\begin{array}{c} Y & D & G & G & A \\ Y & D & G & G & - \\ F & E & G & G & I \\ F & D & - & G & I \\ F & D & - & G & I \\ Y & E & G & G & A \end{array}$	$\begin{array}{c c} characters \\ 1 & 2 & 3 & 4 & 5 & 6 \\ \hline Y & D & G & G & A & V \\ Y & D & G & G & - & - \\ F & E & G & G & I & L \\ F & D & - & G & I & L \\ Y & E & G & G & A & V \\ \end{array}$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c cccc} characters \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \hline Y & D & G & G & A & V & - & E & A \\ Y & D & G & G & - & - & - & E & A \\ F & E & G & G & I & L & V & E & A \\ F & D & - & G & I & L & V & Q & A \\ Y & E & G & G & A & V & V & Q & A \\ \end{array}$

(Can further weight the profile using PAM or BLOSUM matrices)

sequences

Sequence logos



A graphic representation of an aligned set of binding sites. A logo displays the frequencies of bases at each position, as the relative heights of letters, along with the degree of sequence conservation as the total height of a stack of letters, measured in bits of information. Subtle frequencies are not lost in the final product as they would be in a consensus sequence

What can we do with multiple alignments?

- Create (databases of) profiles derived from multiple alignments for protein families
 - profile = multiple alignment + observed character frequencies at each position
- Search with a sequence against a database of profiles (e.g. **PROSITE** database)
 - faster than sequence against sequence
 - gives a more general result ("the input sequence matches globin profile")
- Search with a profile against a database of sequences
 - PSI-BLAST : can identify more distant relationships than by normal BLAST search

PSI-BLAST (position specific iterated BLAST)



PSI-BLAST (Altschul et al 1997)

(1) Start with 1 sequence (or profile) = 'probe'

(2) Search with BLAST and select top hits manually or automatically

(3) Make multiple alignment & profile

(4) Estimate statistical significance of local alignments. If significance ok & you want to continue, then go to (1) using the profile, else **exit**

Dates & programs



Patterns and alternative representations

- Patterns
 - unions of patterns
 - decision trees
 - exact/approximate matching
- Alignments, weight matrices, profiles, HMMs, Neural networks, SCFGA, ...

Brazma et al, Approaches to the automatic discovery of patterns in biosequences, Journal of Computational Biology, 5(2):277-303, 1998

Some terminology

Common similarities between sequences/structures:

- pattern, motif, fingerprint, template, fragment, core, site, alignment, weight matrix, profile...
- "Pattern": description of structure properties
 - (Deterministic) Decide if a protein matches it or not
 - (*Probabilistic*) Assign a value to the match
- "Motif" pattern with biological meaning

Adapted from: Eidhammer, Jonassen & Taylor, "Structure Comparison and Structure Patterns", JCB, 7:5 pp 685-716, 2000.

Classification of functions



Consensus patterns

Alignments

Blocks or Weight Matrices

Templates or Profiles

Bayesian Networks

Hidden Markov Models

Discrete patterns

- Advantages
 - simple and easily interpretable objects
 - easier to discover from scratch (i.e., if no additional information to sequences are given), particularly in noisy data
- Disadvantages
 - limited descriptive power (no weights can be attributed to alternatives)

Regular expressions

- **Symbol**: for each symbol **a** in the alphabet of the language, the regular expression **a** denotes the language containing just the string a
- Alternation: Given 2 regular expressions M and N then M | N is a new regex. A string is in lang(M|N) if it is lang(M) or lang(N). The lang(a|b) = {a,b} contains the 2 strings a and b.
- Concatenation: Given 2 regexes M and N then M•N is a new regex. A string is in lang(M•N) if it is the concatenation of 2 strings α and β s.t. α in lang(M) and β in lang(N). Thus regex (a|b)•a = {aa,ba} defines the language containing the 2 strings aa and ba

Regular expression notation

- a ordinary character, stands for itself
- ε the empty stringanother way to write the empty string!
- $M \mid N$ alternation
- $M \bullet N$ concatenation
- *M*^{*} repetition (zero or more times)
- *M*+ repetition (one or more times)
- *M*? Optional, zero or one occurrence of M
- [a-zA-Z]Character set alternation
- Period stands for any single character except newline
- "a.+*" quotation, string stands for itself

Biosequences - general

- Basic alphabet
 Σ = { a, t/u, c, g} (DNA/RNA)
 Σ = {A, C, ..., Y} (Protein sequence)
- Character group alphabet Π = {g₁...g_n} (e.g. amino-acid class)
- Wild card $X = \{ x (n_1, n_2) | n_1 < n_2 \in N \}$
- $V(x(c_1,c_2))$ set of all words over Σ of length between c_1 and c_2
- Pattern $P = p_1 \dots p_n$, $p_i \in \Sigma \cup \Pi \cup X$

 \rightarrow character & position constraints \leftarrow

Pattern notation and matching

- Separate the pattern alphabet characters by a dash "-"
- Pattern

 $P = A-x(2,6)-[LI]-x(0,\infty)$

matches string

S = ACDEFLGHJKL

because

 $S = A \bullet CDEF \bullet L \bullet GHJKL$

(• meaning concatenation) and

 $A \in V(A)$, $CDEF \in V(x(2,6))$, $L \in V([LI])$, $GHJKL \in V(x(0,\infty))$

PROSITE patterns

- Database of protein families and domains
- Consists of biologically significant sites, patterns and profiles that help to reliably identify to which known protein family (if any) a new sequence belongs
- `x' any amino acid
- Ambiguities :
 - [ALT] =Ala or Leu or Thr
 - **{AM}** any amino acid except Ala and Met.
- `-' separator, `<` N-terminal, `>` C-terminal

٠

- `.` end of pattern
- Repetition: x(3) = x-x-x
- x(2,4) = x-x or x-x-x or x-x-x-x.

PROSITE examples

• $[AC]-x-V-x(4)-\{ED\}.$

- [Ala or Cys]-x-Val-x-x-x-{any but Glu or Asp}

- <A-x-[ST](2)-x(0,1)-V.
 - Start at N-terminal of the sequence
 - Ala-x-[Ser or Thr]-[Ser or Thr]-(x or none)-Val

How to obtain these patterns?

Example property

A given sequence belongs to the chromo-domain family if it matches either the pattern:

```
E-x(0,1)-E-E-[FY]-x-V-E-K-[IV]-[IL]-D-[KR]-R-x(3,4)-G-x-V-
x-Y-x-L-K-W-K-G-[FY]-x-[ED]-x-[HED]-N-T-W-E-P-x(2)-N-
x-[ED]-C-x-[ED]-L-[IL]
```

or the pattern:

L-x(2,3)-E-[KR]-I-[IL]-G-A-[TS]-D-[TSN]-x-G-[EDR]-L-x-F-L-x(2)-[FW]-[KE]-x(2)-D-x-A-[ED]-x-V-x-[AS]-x(2)-A-x(2)-Kx-P-x(2)-[IV]-I-x-F-Y-E

or the pattern:

```
Y-x(0,2)-L-[IV]-K-W-x(6)-[HE]-x-[TS]-W-E-x(4)-[IL]
```

Example family (zinc finger c2h2)



RNA structural patterns

- Constraints:
 - string length
 - inter-string distance
 - character contents
 - matching positions
 - correlation (identical, reverse, complement).
- Complements a-u g-c, g-u (weaker)
- Structures: Stem-loops, Pseudo-knots, Clover leafs
- Context free grammar

Eidhammer, Jonassen, Grinhang, Gilbert & Ratnayake, A contraint-based structure description language for biosequences, Journal of Constraints 6:2/3, 2001

Possible patterns

- Tandem repeat $\alpha \alpha$ <u>acg acg</u>
- Simple repeat $\alpha \beta \alpha$ <u>acgaaaaacg</u>
- Multiple repeat $\alpha \beta \alpha \delta \alpha$ acgaaacguuacg
- Palindrome $\alpha \alpha^r$ <u>acg gca</u>
- Stem loop $\alpha \beta \alpha^{rc}$ <u>acgaacgu</u>
- Pseudoknot $\alpha \gamma_1 \beta \gamma_2 \alpha^{rc} \gamma_3 \beta^{rc}$ augg<u>cuga</u>aggc<u>cgau</u>c<u>ucag</u>ggcau<u>aucg</u>ccgu

(1) $\begin{pmatrix} c \\ a-u \\ g-c \\ u-a \\ c-q \end{pmatrix}$ (2) $\begin{pmatrix} g \\ u-a \\ a-u \\ g-c \\ c-q \end{pmatrix}$

aggc



ggcau

augg

ccgu





(c) David Gilbert 2007

Various ways of using pattern matching for family characterization

A sequence belongs to the family if

- 1. it matches the given sequence *pattern*;
- 2. if it is within a certain *distance* from a string that matches a the pattern (distance between strings can be defined either as a number of mismatches, or as an edit-distance, or based on similarity matrices or some other way);
- 3. if it matches one of a given set of patterns (i.e., if it matches a union of patterns);
- 4. if a decision-tree over the matching patterns returns "yes"

Learning

- Automatically find pattern (given a training set)
- <u>Characterisation</u>: (positive examples only) patterns describing "interesting" properties of a family
- <u>Classification</u>: (positive **and** negative examples) pattern distinguishing S+ and S- .. Which may overlap...

- Formal language for descriptions
- Scoring function to rate descriptions
- Algorithm

Pattern discovery in biosequences

- Motivation:
 - gene functional class prediction
 - RNA splicing
 - protein structure & function
 - gene regulation (transcription factor binding site prediction)
 - detection of repeats

- Prediction of structure/function from sequence:
 - sequence database similarity search
 - compare to family descriptions
 - structure prediction programs

[Alvis Brazma & Inge Jonnassen]

Pattern discovery in biosequences

- Group together sequences thought to have common biological (structural, functional) properties -> families (biological - semantic level)
- Study the purely syntactic properties common to these sequences ignoring their biological (semantic) properties -> patterns, clusters (mathematical - syntactic level)
- Test whether the discovered patterns make sense (back to semantic level)

Protein family analysis

- Collect sequences (structures) in family
- Analyze
 - local multiple alignment
 - global multiple alignment
 - pattern discovery
- Make family description
- Pick up more family members?
 Analyze extended set

Pattern discovery (machine learning)

- Languages & associated discovery mechanisms
- Strings much work
- Finding gene expression sites in DNA may require context sensitive patterns.
- Structures

Approaches to pattern discovery

• Pattern driven:

enumerate all (or some) patterns up to certain complexity (length), for each calculate the score, and report the best

• Sequence driven:

look for patterns by aligning the given sequences

Pattern driven algorithms

- Brute force enumerate all patterns (for instance, all substrings) up to a given length (complexity)
- Evaluate their fitness with respect to the input sequences and output the best
- Unrealistic for patterns of even modest size even for substring patterns (e.g., for substring patterns of length 10 over the amino acid alphabet, there are more than 10¹³ different substrings to enumerate in this way)

Sequence driven algorithms

- Group similar sequences together (e.g., in pairs);
- For each group find a common pattern (e.g., by dynamic programming);
- Group similar patterns together and repeat the previous step until there is only one group left

Sequence driven approach



Algorithm for string pattern discovery

Design (a naive) algorithm for a simple language *s* where s ∈Σ* and * is a wild card of arbitrary length, i.e. x(0,inf)

Example: $s_1 = \text{TAWCEFGOPA}$ $s_2 = \text{FGOPAAWCES}$ $s_3 = \text{WUVTAWCESAW}$

Try discovering patterns using pattern-driven & sequence-driven approaches Sequence-driven: P(s) == set of patterns for s $P(s1) = \{s1\}, P(s2) = \{s2\}, P(s3) = \{s3\}$ $P(s1,s2) = \{...\}, P(s1,s2,s3) = \{...\}$

Amino acid residue groups

Residue property

Residue groups

Small Small hydroxyl Basic Aromatic Basic+ Small hydrophobic Medium hydrophobic Acidic/amide Small/polar

Ala, Gly	A,G
Ser, Thr	S,T
Lys, Arg	K,R
Phe, Tyr, Trp	F,Y,W
His, Lys, Arg	H,K,R
Val, Leu, Ile	V,L,I
Val, Leu, Ile, Met	V,L,I,M
Asp, Glu, Asn, Gln	D,E,N,Q
Ala, Gly, Ser, Thr, Pro	A,G,S,T,P

Deriving regular expressions

$$\begin{split} s_1 &= \texttt{ALDGAVFALCDRYFQ} \\ s_2 &= \texttt{SDVGPRSCFCERFYQ} \\ s_3 &= \texttt{ADLGRTQNRCDRYYQ} \\ s_4 &= \texttt{ADIGQPHSLCERYFQ} \end{split}$$

Make a regular expression & a 'fuzzy' regular expression!



Rating patterns

- Size (e.g. number of characters...).
 - Hence Information content: e.g. length of the pattern (& perhaps penalties for wild cards)
- Compression
 - measure of how much of each of the items in the learning set is described
- Sensitivity, Specificity etc
 - requires evaluation against learning [training] & test sets

Compression - see updated slides

Send the pattern once, and then for each item, send the unmatched parts

(1) Raw Compression (chars k):

$$C_{raw} = (\sum_{i \in 1..n} N(k_i)) - (n-1)*N(k_p)$$

sum of chars in the examples minus (No_examples - 1) * chars_in_pattern Varies from ? to ?

(2) Normalised compression: $C_{norm} = 1 - ((\sum_{i \in 1..n} N(k_i)) - C_{raw}) / ((\sum_{i \in 1..n} N(k_i)) - min(N(k_i)))$

This is a goodness of compression measure (0=good to 1=bad).

(c) David Gilbert 2007

Compression

Send the pattern once, and then for each item, send the unmatched parts

1

$$C_{raw} = \left(\sum_{i=1}^{n} |S_i|\right) - (n-1) * |P$$

 \mathbf{N}

i.e. SumOfElementsInExamples - (NumberOfExamples - 1) * elements in pattern

(2) Normalised compression:

(1) *Raw Compression*:

$$C_{norm} = \frac{\left(\sum_{i=1}^{n} |S_i|\right) - C_{raw}}{\left(\sum_{i=1}^{n} |S_i|\right) - \min_{i=1}^{n} \left(|S_i|\right)}$$

This is a goodness measure (1=good, 0=bad). (c) David Gilbert 2007 Multiple Alignment, Patterns & Profiles

More compression

(3) Substituting (1) into (2):



(4) Pairwise comparison via compression:



Characteristic string function for family F+



Classification & conservation problems



Classification problem C1

• Given a set S+ of sequences believed to be members of family F+, and a set S- of sequences believed not to be members, i.e.

 $S+ \subset F+ and S- \subset F-$

 $F+ \cap F- = \emptyset$ and $F+ \cup F- = \Sigma^*$

- Find *compact* string functions that return
 - TRUE for all $s \in S+$ and FALSE for all $s \in S-$, and
 - have a high likelihood for returning TRUE for $s \in F+$ and FALSE for $s \in F-$
- C1a: find compact "explanations" of known sequences
- C1b: try to predict the family relationship of yet unknown sequences
- N1: suppose $F^+ \cap F^- = \emptyset$ and $F^+ \cup F^- = \Sigma^* i$ and $S^+ \cap F^-$ and $S^- \cap F^+$ are small, find *compact* string functions that return
 - TRUE for *most* $s \in S^+$ and FALSE for *most* $s \in S^-$, and
 - have a high likelihood for returning TRUE for $s \in F+$ and FALSE for $s \in F-$

Characterisation: conservation problem C2

- Given a set S+ of sequences believed to be members of family F+, i.e. S+ \subset F+
- Find *interesting* string functions that return
 - TRUE for all $s \in S^+$
 - have a high likelihood for returning TRUE for $s \in F+$
- N2: suppose $F+ \subset \Sigma^*$, and given $S+ \subset \Sigma^*$, such that $S+ \cap (F+)$ is small, find *interesting* string functions that return
 - TRUE for *most* $s \in S^+$, and
 - have a high likelihood for returning TRUE for $s \in F^+$
- *Interesting*: have a low probability for returning TRUE for random sequences

Training and test sets

• *training* set of

S+ positive examples from F+, and optionally a set S- of negative examples from F-

• *test* set

T+ from F+ where T+ \cap S+ = \emptyset , and optionally T- from F- where T- \cap S- = \emptyset

- In practice, we may not know all members of F+ and F-
 - Thus to construct training & test sets, we can randomly divide an initial set of positive examples into a training set S+ and a test set T+, similarly for S- and T-
 - The goal is to accurately describe "new" members of F+ and F- when we come across them

Training and test sets





The challenge of increasing data "All data" Training Set Current data (continues to expand) Language of the pattern L(P)

True positives, true negatives, false positives, false negatives





F-measure = 2 * (Precision * Recall) / (Precision + Recall)

F-measure

 F_1 -measure = 2 * (Precision * Recall) / (Precision + Recall)

General F-measure = $(1+\alpha)$ * (Precision * Recall) / (α *Precision + Recall)

Training and test sets (positive examples only)



Methodology

- Solution space / hypothesis space / target class: find a good class of string functions from which the approximating function *f* is chosen for a real-world problem
- *Fitness measure*: define a ranking of the solution space, evaluating how good each function is for the training set (how likely *f* is to approximate *g*
- Develop an *algorithm* returning those classifier functions from the given solution space that rate high enough according to the fitness measure

Defining string functions via patterns

Given a string *s* and a pattern π which defines a language $L(\pi)$, define a classification (conservation) function *f* by

$$f(s) = \begin{cases} \text{TRUE if } s \in L(\pi) \\ \text{FALSE otherwise} \end{cases}$$

$$f(s) = \begin{cases} \text{TRUE if } Dist(\pi, s) \leq const \\ \text{FALSE otherwise} \\ \text{Where } Dist(\pi, s) = \min_{s' \in L(\pi)} dist(s', s) \end{cases}$$

Clean / Noisy Data

- <u>Clean data</u>: the training set is assumed to be "correct"
- <u>Noisy data</u>: training set
 - sequences may contain errors
 - sequences may have been assigned to the wrong family

PROSITE profiles

• Uses Hidden Markov Model - can characterise an entire family of sequences.

