

Getting To Know Your Protein

Comparative Protein Analysis:

Part II. Protein Domain Identification & Classification

Robert Latek, PhD Sr. Bioinformatics Scientist Whitehead Institute for Biomedical Research

Syllabus

• Protein Families

- Identifying Protein Domains
- Family Databases & Searches
- Searching for Family Members
 - Pattern Searches
 - Patscan
 - Profile Searches
 - PSI-BLAST/HMMER2

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

Proteins As Modules

Comparative Protein Analysis

to understand global relationships between sequences.

Phylogenetic Trees and Multiple Sequence Alignments are important tools

http://jura.wi.mit.edu/bio/education/bioinfo-mini/seq/ (AA Subst. Matrices)
 Part II.:
 How do you identify sequence relationships that are restricted to localized

Can you apply phylogenetic trees and MSAs to only sub-regions of

– How do you apply what you know about a group of sequences to finding additional, related sequences?

What can the relationship between your sequences and previously discovered ones tell you about their function?

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

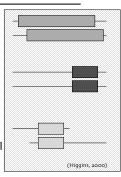
Part I.:

sequences'

 Proteins are derived from a limited number of basic building blocks (Domains)

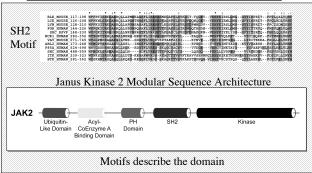
Assigning sequences to Protein Families

- Evolution has shuffled these modules giving rise to a diverse repertoire of protein sequences
- As a result, proteins can share a global or local relationship



WIBR Bioinformatics Courses, © Whitehead Institute, 2004

Protein Domains



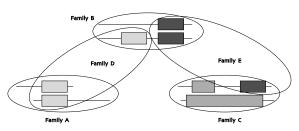
WIBR Bioinformatics Courses, © Whitehead Institute, 2004

Protein Families

- Protein Family a group of proteins that share a common function and/or structure, that are potentially derived from a common ancestor (set of homologous proteins)
- Characterizing a Family Compare the sequence and structure patterns of the family members to reveal shared characteristics that potentially describe common biological properties
- Motif/Domain sequence and/or structure patterns common to protein family members (a trait)

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

Protein Families

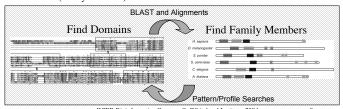


Separate Families can Be Interrelated

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

Creating Protein Families

- Use domains to identify family members
 - Use a sequence to search a database and characterize a pattern/profile
 - Use a specific pattern/profile to identify homologous sequences (family members)



WIBR Bioinformatics Courses, © Whitehead Institute, 2004

Family Database Resources

- Curated Databases*
 - Proteins are placed into families with which they share a specific sequence pattern
- Clustering Databases*
 - Sequence similarity-based without the prior knowledge of specific patterns
- Derived Databases*
 - Pool other databases into one central resource
- · Search and Browse

*(Higgins, 2000)

11

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

Curated Family Databases • Pfam (http://pfam.wustl.edu/hmmsearch.shtml/) **

- - Uses manually constructed seed alignments and PSSM to automatically extract domains
 - db of protein families and corresponding profile-HMMs of prototypic domains
 - Searches report e-value and bits score
- Prosite (http://www.expasy.ch/tools/scanprosite/)
 - Hit or Miss -> no stats
- PRINTS (http://www.bioinf.man.ac.uk/fingerPRINTScan/)



WIBR Bioinformatics Courses, © Whitehead Institute, 2004

Clustering Family Databases

- Search a database against itself and cluster similar sequences into families
- ProDom (http://prodes.toulouse.inra.fr/prodom/current/html/home.php)
 - Searchable against MSAs and consensus sequences
- Protomap (http://protomap.cornell.edu/)

Swiss-Prot based and provides a tree-like view of clustering



WIBR Bioinformatics Courses, © Whitehead Institute, 2004

Derived Family Databases

- Databases that utilize protein family groupings provided by other resources
- Blocks Search and Make (http://blocks.fhcrc.org/blocks/)
 - Uses Protomap system for finding blocks that are indicative of a protein family (GIBBS/MOTIF)
- Proclass (http://pir.georgetown.edu/gfserver/proclass.html)
 - Combines families from several resources using a neural network-based system (relationships)
- MEME (http://meme.sdsc.edu/meme/website/intro.html)



Syllabus

- Protein Families
 - Identifying Protein Domains
 - Family Databases & Searches
- Searching for Family Members
 - Pattern Searches
 - Patscan
 - Profile Searches
 - PSI-BLAST/HMMER2

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

13

Searching Family Databases

- BLAST searches provide a great deal of information, but it is difficult to select out the important sequences (listed by score, not family)
- Family searches can give an immediate indication of a protein's classification/function
- Use Family Database search tools to identify domains and family members

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

Patterns & Profiles

- Techniques for searching sequence databases to uncover common domains/motifs of biological significance that categorize a protein into a family
- **Pattern** a deterministic syntax that describes multiple combinations of possible residues within a protein string
- Profile probabilistic generalizations that assign to every segment position, a probability that each of the 20 aa will occur

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

Discovery Algorithms

- Pattern Driven Methods
 - Enumerate all possible patterns in solution space and try matching them to a set of sequences

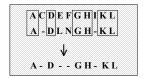


WIBR Bioinformatics Courses, © Whitehead Institute, 2004

16

Discovery Algorithms

- Sequence Driven Methods
 - Build up a pattern by pair-wise comparisons of input sequences, storing positions in common, removing positions that are different



Pattern Building

- Find patterns like "pos1 xx pos2 xxxx pos3"
 - Definition of a non-contiguous motif



Define/Search A Motif http://us.expasy.org/tools/scanprosite/

Pattern Properties

Specification

a single residue K, set of residues (KPR), exclusion {KPR}, wildcards X, varying lengths x(3,6) -> variable gap lengths

General Syntax

- C-x(2,4)-C-x(3)-[LIVMFYWC]-x(8)-H-x(3,5)-H
- Patscan Syntax (http://jura.wi.mit.edu/bio/education/bioinfo-mini/seq)
 - C 2...4 C 3...3 any(LIVMFYWC) 8...8 H 3...5 H

• Pattern Database Searching

- %scan_for_matches -p pattern_file < nr > output_file

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

10

Sequence Pattern Concerns

- Pattern descriptors must allow for approximate matching by defining an acceptable distance between a pattern and a potential hit
 - Weigh the sensitivity and specificity of a pattern
- What is the likelihood that a pattern would randomly occur?

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

20

Sequence Profiles

- Consensus mathematical probability that a particular aa will be located at a given position
- Probabilistic pattern constructed from a MSA
- Opportunity to assign penalties for insertions and deletions, but not well suited for variable gap lengths
- **PSSM** (Position Specific Scoring Matrix)
 - Represents the sequence profile in tabular form
 - Columns of weights for every aa corresponding to each column of a MSA

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

21

Profile Analysis

- · Perform global MSA on group of sequences
- · Move highly conserved regions to smaller MSAs
- Generate scoring table with log odds scores
 - Each column is independent
 - Average Method: profile matrix values are weighted by the proportion of each amino acid in each column of MSA
 - Evolutionary Method: calculate the evolutionary distance (Dayhoff model) required to generate the observed amino acid distribution



WIBR Bioinformatics Courses, © Whitehead Institute, 2004

PSSM Example V T M G (i.e. Distribution of aa in an MSA column) Target sequences Resulting Consensus: ITLS PSSM



WIBR Bioinformatics Courses, © Whitehead Institute, 2004

PSSM Properties

- Score-based sequence representations for searching databases
- Goal
 - Limit the diversity in each column to improve reliability
- Problems
 - Differing length gaps between conserved positions (unlike patterns)

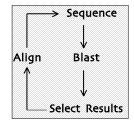
WIBR Bioinformatics Courses, © Whitehead Institute, 2004

PSI-BLAST Implementation

• PSI-BLAST

http://www.ncbi.nlm.nih.gov/BLAST/

- Start with a sequence, BLAST it, align select results to query sequence, estimate a profile with the MSA, search DB with the profile - constructs PSSM
- Iterate until process stabilizes
- Focus on domains, not entire sequences
- Greatly improves sensitivity



WIBR Bioinformatics Courses, © Whitehead Institute, 2004

25

PSI-BLAST Sample Output

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

_ .

HMM Building

- Hidden Markov Models are Statistical methods that consider all the possible combinations of matches, mismatches, and gaps to generate a consensus (Higgins, 2000)
- Sequence ordering and alignments are not necessary at the onset (but in many cases alignments <u>are</u> recommended)
- Ideally use at least 20 sequences in the training set to build a model
- Calibration prevents over-fitting training set (i.e. Ala scan)
- Generate a model (profile/PSSM), then search a database with it

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

27

HMM Implementation

- HMMER2 (http://hmmer.wustl.edu/)
 - Determine which sequences to include/exclude
 - Perform alignment, select domain, excise ends, manually refine MSA (pre-aligned sequences better)
 - Build profile
 - | %hmmbuild [-options] <hmmfile output> <alignment file>
 - Calibrate profile (re-calc. Parameters by making a random db)
 - %hmmcalibrate [-options] <hmmfile>
 - Search database
 - %hmmsearch [-options] <hmmfile> <database file> > out

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

2

HMMER2 Output

- Hmmsearch returns evalues and bits scores
- Repeat process with selected results
 - Unfortunately need to extract sequences from the results and manually perform MSA before beginning next round of iteration



Patterns vs. Profiles

• Patterns

- Easy to understand (human-readable)
- Account for different length gaps
- Profiles
 - Sensitivity, better signal to noise ratio
 - Teachable

Domain ID & Searching

- Family/Domain Search
 - http://pfam.wustl.edu



- · Pattern Search
 - scan_for_matches (Patscan)
 - scan_for_matches -p pattern_file </cluster/db0/Data/yeast.aa > output_file
 any(CF) any(HF) any(GK) 1...1 any(LJ) 4...4 any(AS) 3...3 any(LJ) 3...3 any(GA) 3...3 G 1...1 any(FF) 1...1 any(LJ) R
- · Profile Search
 - HMMER2
 - hmmbuild [-options] <hmmfile output> <alignment file>

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

33

References

- Bioinformatics: Sequence and genome Analysis. David W. Mount. CSHL Press, 2001.
- Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins. Andreas D. Baxevanis and B.F. Francis Ouellete. Wiley Interscience, 2001.
- Bioinformatics: Sequence, structure, and databanks. Des Higgins and Willie Taylor. Oxford University Press, 2000.

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

Exercises

- Use PFAM to identify domains within your sequence
- Scan your sequences with ProSite to find a pattern to represent the domain
- Use the ProSite pattern to search the non-redundant db
- Use PSI-BLAST to build a sequence profile and search the non-redundant db

WIBR Bioinformatics Courses, © Whitehead Institute, 2004

32