



Sequence Analysis

II: Sequence Patterns and Matrices

George Bell, Ph.D.
WIBR Bioinformatics and Research Computing

Sequence Patterns and Matrices

- Multiple sequence alignments
- Sequence patterns
- Sequence matrices
- Identifying regulatory sites
- Finding overrepresented patterns and profiles
- Gene finding

Sequence Analysis Course © Whitehead Institute, 2005

2

Why use DNA patterns and matrices?

- To help search the genome for ...
 - Transcription start sites
 - Splice junctions (exon-intron boundaries)
 - Transcription factor binding sites
 - microRNA targets
 - Active sites for chromatin regulators
 - Gene regions encoding protein motifs
 - RNA folding patterns (hairpins, etc.)

Sequence Analysis Course © Whitehead Institute, 2005

3

Multiple sequence alignments (MSAs)

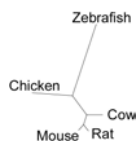
- Global MSA is computationally difficult
- As a result, MSA algorithms use approximate methods
- Independent of the chosen algorithm, choice of scoring matrix is important
- Aligning contigs vs. genes
- Aligning similar vs. divergent sequences

Sequence Analysis Course © Whitehead Institute, 2005

4

Global progressive MSA

- An MSA method that uses phylogenetic information to determine alignment order
1. Perform all pairwise alignments
 2. Use alignment scores to create a tree
 3. Align most similar pair of sequences and create consensus.
 4. Align next most similar pair of sequences and create consensus ... repeat until done



```

Rat  GAATGTTGGATTGGCC
Mouse GAATGTTGGATTGGCC
Cow  GAATGTTGGATTGGCC
Chicken GAATGTTGGATTGGCC
Zebrafish GAATGTTGGATTGGCC
  
```

Sequence Analysis Course © Whitehead Institute, 2005

5

Sequence patterns

Pattern: an expression describing all possible combinations of bases in a sequence

- Generally derived from a MSA
- Ex1. EcoRI enzyme site: GAATTC
- Ex2. Codons for proline: CC[ACTG]; CCN
- Ex3. TATA box: TATA[AT][AGT][GA]
- Ex4. TFBS for GATA4: AGATA[AGT][AC]AGGGA
- Ex 5. Gene region encoding your favorite protein motif => better to use protein pattern!

Sequence Analysis Course © Whitehead Institute, 2005

6

More complex patterns

- May want to consider:
 - Mismatches
 - Insertions
 - Deletions
 - Alphabet reflecting ambiguity
- Ex: Patscan (Argonne National Laboratory) syntax
 - Pattern[Mismatches,Deletions,Insertions]
 - Ex: RRRRRYYYYY[3,2,1]
(R = purine; Y = pyrimidine)

Pattern considerations

- Is there reliable data behind it?
- Is it specific and sensitive?
- How many matches would you expect by chance?
- Patterns don't represent the different probabilities of each combination of bases, just whether they can occur or not.
- DNA or protein?

Pattern searching programs

- Check examples or help for syntax
- EMBOSS:
 - fuzznuc: nucleic acid pattern search
 - fuzzpro: protein pattern search
 - dreg: regular expression search of a nucleotide sequence
- PatScan
- Perl (programming language) regular expressions

Sensitivity and Specificity

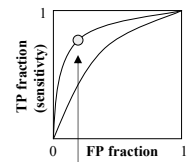
Proportion of true sites correctly identified:

$$\text{sensitivity} = \frac{TP}{TP + FN}$$

Proportion of false sites correctly identified:

$$\text{specificity} = \frac{TN}{TN + FP}$$

ROC plot:



Aim for "optimal" sensitivity

Matrix Representations

Matrix: a probabilistic representation of bases in a sequence

- Generally derived from a MSA
- Related to concept of "profile" (but no gaps allowed in MSA)
- Maintains meaning when transposed
- Position-specific scoring matrix (PSSM) assumes each position is independent
- Handling "zero" probabilities with pseudocounts

Creating a matrix (PSSM)

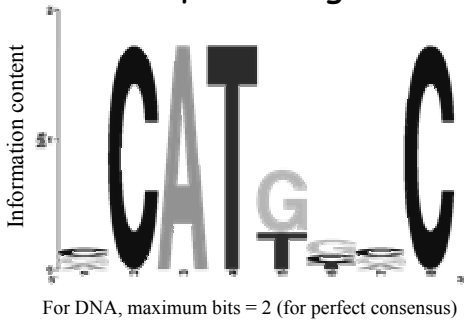
1. Create alignment
2. Count frequencies; add pseudocounts

A	GCATTGCG		1	2	3	4	5	6	7	8	$\psi = \frac{\sqrt{n_{\text{seq}}}}{n_{\text{seq}}}$ $= \frac{\sqrt{8}}{17}$ $= 0.167$	
B	ACATGGAC		A	2	ψ	8	ψ	ψ	ψ	2		ψ
C	CCATGGCC		C	3	8	ψ	ψ	ψ	2	3		8
D	ACATGGAC		G	5	ψ	ψ	ψ	5	4	3		ψ
E	CCATTGCC		T	ψ	ψ	ψ	8	3	2	ψ		ψ

3. Calculate log-odds scores: $\log_2(\text{freq}_{\text{obs}}/\text{freq}_{\text{exp}})$

	1	2	3	4	5	6	7	8
A	-2.6	-2.6	2	-2.6	-2.6	-2.6	-2.6	-2.6
C	0.6	2	-2.6	-2.6	-2.6	-2.6	0.6	2
G	1.3	-2.6	-2.6	-2.6	1.3	1	0.6	-2.6
T	-2.6	-2.6	-2.6	2	0.6	-2.6	-2.6	-2.6

Sequence logos



Sequence Analysis Course © Whitehead Institute, 2005

13

Searching with matrices

- Slide matrix along sequence(s), take sum of log-odds scores for base at each sequence position:

Example with sample matrix of length 3, showing scores at two positions

	1	2	3										
A	-2.6	-2.6	2		A	-2.6	-2.6	2					
C	0.6	2	-2.6		C	0.6	2	-2.6					
G	1.3	-2.6	-2.6		G	1.3	-2.6	-2.6					
T	-2.6	-2.6	-2.6		T	-2.6	-2.6	-2.6					
G	T	A	C	G	A	C	G	T	G	C	C	A	T

$$\Sigma = (-2.6) + 2 + (-2.6) = -3.2$$

$$\Sigma = 0.6 + 2 + 2 = 4.6$$

Highest score wins

Sequence Analysis Course © Whitehead Institute, 2005

14

Hints for identifying regulatory sites

- Mask repetitive sequence first (RepeatMasker) to remove “non-functional noise”
- What specific area(s) of the sequence or genome can you search (instead of all of it)?
- Look at conservation: functional regulatory sites tend to be conserved
- ENCODE project (1% of human genome)

Sequence Analysis Course © Whitehead Institute, 2005

15

Identifying over-represented patterns

- Count oligos of each sequence of expected length.
- Calculate expected frequencies.
- Rank observed/expected values.
- Repeat for oligos of another length.

This method assumes the pattern is very specific

Sequence Analysis Course © Whitehead Institute, 2005

16

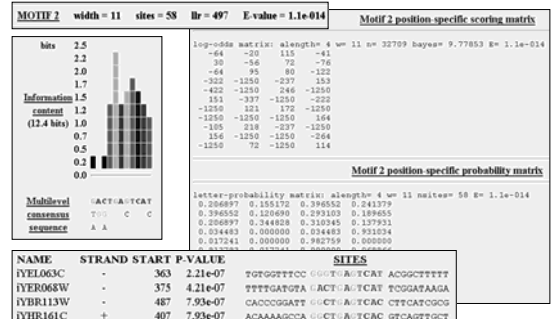
Identifying over-represented matrices

- Inputs
 - a set of sequences assumed to contain a matrix
 - range of presumed profile width?
 - ≥ 0 (“zoops”) or ≥ 1 (“oops”) occurrence per sequence?
- Programs
 - Meme: based on the expectation maximum (EM) algorithm; meme.sdsc.edu
 - AlignACE: based on the Gibbs sampling algorithm; atlas.med.harvard.edu

Sequence Analysis Course © Whitehead Institute, 2005

17

Meme sample output



Sequence Analysis Course © Whitehead Institute, 2005

18

Identifying features of genes in genomic DNA

- Splice sites
- Open reading frames
- Promoters
- Codon bias
- Expression information (ESTs, mRNA)
- Protein similarity to known genes
- Conservation across species

Gene finding programs (sample)

- GeneWise (Birney and Durbin, 2000)
 - Genscan (Burge and Karlin, 1997)
 - Acembly (Thierry-Mieg et al.)
 - Twinscan (Korf et al., 2001)
 - SGP (Parra et al., 2003)
 - GeneID (Parra et al., 2000)
- Use all available data and predictions when possible

Summary

- Multiple sequence alignments
- Sequence patterns
- Sequence matrices
- Identifying regulatory sites
- Finding over-represented patterns and matrices
- Gene finding

References

- Bioinformatics: Sequence and Genome Analysis, 2nd ed. David Mount. CSHL Press, 2004.
- Publications describing algorithms and software for
 - multiple sequence alignment
 - pattern and matrix analysis and searching
 - gene finding

Exercises

1. Investigating the mechanisms of miRNA activity through pattern searching
2. Studying transcriptional control with DNA matrices

Both involve computational analysis of data from recently published studies