For written notes on this lecture, please read chapter 10 of The Practical Bioinformatician

CS2220: Introduction to Computational Biology Lecture 6: Essence of Sequence Comparison

> Limsoon Wong 16 February 2007







- Dynamic Programming
- String Comparison
- Sequence Alignment
 - Pairwise Alignment
 - Needleman-Wunsch global alignment algorithm
 - Smith-Waterman local alignment algorithm
 - Multiple Alignment
- Popular tools
 - FASTA, BLAST, Pattern Hunter
- More advanced tools

- PHI BLAST, ISS, PSI-BLAST, SAM, ...

What is Dynamic Programming



NUS National University of Singapore

The Knapsack Problem

Source: http://mat.gsia.cmu.edu/classes/dynamic/node6.html

- The problem
 - Each item that can go into the knapsack has a size and a benefit. The knapsack has a certain capacity. What should go into the knapsack so as to maximize the total benefit?
- A dynamic programming solution
 - Let w_j and b_j be weight and benefit for item *j*. Let g(w) be max benefit that can be gained from a *w* pound knapsack. Then g(w) relates to previously calculated *g* values as follows:

$$g(w) = \max_{j} \{b_j + g(w - w_j)\}$$

An Example



Source: http://mat.gsia.cmu.edu/classes/dynamic/node6.html

• Suppose the items are

Item (j)	Weight (wy)	Benefit(by)
1	2	66
2	8	80
8	1	80

• Recall that

$$g(w) = \max_j \{b_j + g(w - w_j)\}$$

 To fill a *w* pound knapsack, we must end off by adding some item. If we add item *j*, we end up with a knapsack of size *w* – *wj* to fill ... • To illustrate:

- g(0) = 0
- g(1) = 30, item 3
- $\begin{array}{l} \quad g(2) = \max\{65 + g(0) = 65, \, 30 + g(1) \\ = 60\} = 65, \, \text{item 1} \end{array}$
- $\begin{array}{l} \quad g(3) = \max\{65 + g(1) = 95, \, 80 + g(0) \\ = 80, \, 30 + g(2) = 95\} = 95, \, \mbox{item 1/3} \end{array}$
- $g(4) = max\{65 + g(2) = 130, 80 + g(1) = 110, 30 + g(3) = 125\} = 130,$ item 1
- $\begin{array}{l} & g(5) = \max\{65 + g(3) = 160, \, 80 + \\ g(2) = 145, \, 30 + g(4) = 160\} = 160, \\ & \textbf{item 1/3} \end{array}$
- ⇒ For knapsack of capacity 5, max benefit is 160, which is gained by adding 2 of item 1 and 1 of item 3

g(1), g(2), ... are computed many times of Singapore



"Memoize" to avoid recomputation





Non-Recursive Version



Characteristics of Dynamic Programming

Source: http://mat.gsia.cmu.edu/classes/dynamic/node4.html

- The problem can be divided into stages with a decision required at each stage Exercise: What is a stage
- Each stage has a number of states associated
 in the Knapsack problem?
- The decision at one stage transforms one state into a state in the next stage in the Knapsack problem?
- Given current state, the optimal decision for each remaining states does not depend on previous states or decisions
 E.g., g(2) doesnt depends on g(3)
- There is a recursive relationship that identifies the optimal decision for stage *j*, given stage *j*+1 has already been solved
- The final stage must be solvable by itself

Copyright 2007 © Limsoon Wong

E.g., g(0) = 0

Sequence Alignment



Motivations for Sequence Comparison Singapore

- DNA is blue print for living organisms
- \Rightarrow Evolution is related to changes in DNA
- ⇒ By comparing DNA sequences we can infer evolutionary relationships between the sequences w/o knowledge of the evolutionary events themselves
- Foundation for inferring function, active site, and key mutations



 Doolittle et al. (Science, July 1983) searched for platelet-derived growth factor (PDGF) in his own DB. He found that PDGF is similar to v-sis oncogene

PDGF-21SLGSLTIAEPAMIAECKTREEVFCICRRL?DR??34p28sis61LARGKRSLGSLSVAEPAMIAECKTRTEVFEISRRLIDRTN100



Sequence Alignment



Sequence Alignment: Poor Example

Poor seq alignment shows few matched positions
 The two proteins are not likely to be homologous

Alignment by FASTA of the sequences of amloyanin and domain 1 of ascorbate oxidase

60 70 80 100Amicyanin **MPHNVHFVAGVLGEAALKGPMMKKEOAYSLTFTEAGTYDYHCTPHPFMRGKVVVE** 11. 11 Ascorbate Oxidase ILQRGT9WADGTASISQCAINPGETFFYNFTVDNPGTFFYHGHLGMQRSAGLYGSLI 7080 90 100 110120 No obvious match between Amicyanin and Ascorbate Oxidase

Sequence Alignment: Good Example Singapore

- Good alignment usually has clusters of extensive matched positions
- \Rightarrow The two proteins are likely to be homologous

D >gil13476732|ref|NP_108301.1| unknown protein [Mesorhizobium loti]
gil14027493|dbj|BAB53762.1| unknown protein [Mesorhizobium loti]
Length = 105

```
Score = 105 bits (262), Expect = 1e-22
Identities = 61/106 (57%), Positives = 73/106 (68%), Gaps = 1/106 (0%)
```

Query: 1 MKPGRLASIALAIIFLPMAVPAHAATIEITMENLVISPTEVSAKVGDTIRWVNKDVFAHT 60 MK G L ++ MA PA AATIE+T++ LV SP V AKVGDTI WVN DV AHT Sbjct: 1 MKAGALIRLSWLAALALMAAPAAAATIEVTIDKLVFSPATVEAKVGDTIEWVNNDVVAHT 60

> good match between Amicyanin and unknown M. loti protein

Alignment: Simple-Minded Probability & Score

Let p, q, r be respectively the probability of a match, a mismatch, and an indel. Then the probability of an alignment A = (X, Y) is

$$prob(A) = p^m \cdot q^n \cdot r^h$$

where

$$\begin{array}{lll} m & = & |\{i \mid x'_i = y'_i \neq -\}| \\ n & = & |\{i \mid x'_i \neq y'_i, x'_i \neq -, y'_i \neq -\}| \\ h & = & |\{i \mid x'_i = -, y'_i \neq -\} \cup \{i \mid x'_i \neq -, y'_i = -\}| \end{array}$$

• Define score S(A) by simple log likelihood as

- S(A) = log(prob(A)) [m log(s) + h log(s)], with log(p/s) = 1
- Then S(A) = #matches μ #mismatches δ #indels

Exercise: Derive μ and δ

Global Pairwise Alignment: Problem Definition



• Given sequences *U* and *V* of lengths *n* and *m*, then number of possible alignments is given by

- f(n, m) = f(n-1,m) + f(n-1,m-1) + f(n,m-1)

 $-f(n,n) \sim (1 + \sqrt{2})^{2n+1} n^{-1/2}$

Exercise: Explain the recurrence above

• The problem of finding a global pairwise alignment is to find an alignment A so that S(A) is max among exponential number of possible alternatives

Global Pairwise Alignment: Dynamic Programming Solution



- s(x,x) = 2
- $s(x,y) = -\mu$, if $x \neq y$
- Then

Let U and V be two sequences of length n and m. Then their global pairwise alignment can be extracted from the dynamic programming computation of $S_{n,m}$, where

$$S_{i,j} = \max \left\{ \begin{array}{c} S_{i-1,j-1} + s(u'_i, v'_j) \\ S_{i-1,j} - \delta \\ S_{i,j-1} - \delta \end{array} \right\}$$

Exercise: What is the effect of a large δ ?

This is the basic idea of the Needleman-Wunsch algorithm





- Source: Ken Sung
- Consider two strings S[1..n] and T[1..m]
- Let V(i, j) be score of opt alignment betw S[1..i] and T[1..j]
- Basis:
 - -V(0, 0) = 0
 - $V(0, j) = V(0, j 1) \delta$
 - Insert j times
 - $-V(i, 0) = V(i 1, 0) \delta$
 - Delete i times

Needleman-Wunsch Algorithm (II)

Recurrence: For i>0, j>0

- $V(i, j) = \max \begin{cases} V(i-1, j-1) + s(S[i], T[j]) & \text{Match/mismatch} \\ V(i-1, j) \delta & \text{Delete} \\ V(i, j-1) \delta & \text{Insert} \end{cases}$
- In the alignment, the last pair must be either match/mismatch, delete, insert





Examp	le (I)
Source: Ken Sun	g

	_	Α	G	С	Α	Т	G	С
_	0	-1	- 2	- 3	- 4	- 5	- 6	- 7
Α	- 1							
С	- 2							
Α	- 3							
Α	- 4							
т	- 5							
С	- 6							
С	- 7							



Example (II) Source: Ken Sung									
	_	Α	G	С	Α	Т	G	С	
_	0	_ –1 ₊	_ - 2 _∢	3_	_ – 4	_ – 5₊	6,	7	
Α	- 1	2							
C.	-	$\int S_{c}$),0 +	- s(2	4, <i>A</i>)		0	÷	2
S ₁ ,		$ax \{ S_0 \}$	D,1		1 -	= max	x - 1		1 = 2
	<u> </u>	S_1	,0 —		1		[-1		
Т	- 5								
С	- 6								
С	- 7								



E	Xam Source: Ke	ple In Sung	()	

	_	Α	G	С	Α	Т	G	С	
_	0	1 _	_ - 2 _∢	3_	4_	_ – 5,_	_ - 6 _∢	7	
Α	- 1	2	_ 1						
internet in the second se		$\int S_0$),1 +	s(A	(,G)		[-1	+ -	-1
S _{1,2}	$_2 = m$	$ax \{ S_0 \}$,2 —		1 =	= max	$\{-2$		1 =1
A.	4	$\int S_1$,1				2		1
Т	- 5								
С	- <mark>6</mark>								
С	- 7								



Exam	ple	(IV)
Source: Ke	en Sung		· · · ·

	_	Α	G	С	Α	т	G	С		
_	0	- 1	- 2	- 3	- 4	- 5	- 6	- 7		
Α	- 1	2	1	0	- 1	- 2	- 3	- 4		
С	- 2	1	1	3	2					
Α	- 3									
Α	- 4									
Т	- 5									
С	- 6	Exercise: Can you tell from these entries what Are the values of $s(A,G)$, $s(A,C)$, $s(A,A)$, etc.?								
С	- 7									



Exam	ple	(V)
Source: Ken	Sung	

	_	Α	G	С	Α	т	G	С
_	0	_ −1 ₊	_ - 2 _∢	3_	_ – 4	_ – 5₊	6₊	7
Α	- 1	2	_ 1 🗸	0	_ – 1,	_ – 2	3,	4
С	- 2	1	1	3	2	_ 1 _	0	1
Α	-3	0	0	2	5 ←	- 4 +	- 3 +	- 2
Α	-4	- 1	- 1	1	4	4	3	2
Т	−5	-2	- 2	0	3	6 ← ↑ `	- 5 +	- 4
С	6	-3	- 3	0	2	5	5	7
С	- 7	- 4	- 4	-1	1	4	4	7





```
Create the table V[0...,0...m] and P[1...,1...m];
V[0,0] = 0;
For j=1 to m, set V[0,j] := v[0,j-1] - \delta;
For i=1 to n, set V[i,0] := V[i - 1,0] - \delta;
For j=1 to m {
  For i = 1 to n {
       set V[i,i] := V[i,i - 1] - \delta ;
       set P[i,j] := (0, -1);
       if V[i,i] < V[i - 1,i] - \delta then
              set V[i,i] := V[i - 1,i] - \delta;
              set P[i, j] := (-1, 0);
       if (V[i,j] < V[i - 1, j - 1] + s(S[i],T[j])) then
              set V[i,j] := V[i - 1, j - 1] + s(S[i],T[j]);
              set P[i,j] := (-1, -1);
}
Backtracking P[n,m] to P[0,0] to find optimal alignment;
```





- We need to fill in all entries in the n×m matrix
- Each entry can be computed in O(1) time
- ⇒ Time complexity = O(nm)
- ⇒ Space complexity = O(nm)

Exercise: Write down the memoized version of Needleman-Wunsch. What is its time/space complexity?

Problem on Speed

Source: Ken Sung

- Aho, Hirschberg, Ullman 1976
 - If we can only compare whether two symbols are equal or not, the string alignment problem can be solved in Ω(nm) time

• Hirschberg 1978

 If symbols are ordered and can be compared, the string alignment problem can be solved in Ω(n log n) time

- Masek and Paterson 1980
 - Based on Four-Russian's paradigm, the string alignment problem can be solved in O(nm/log2 n) time
- Let d be the total number of inserts and deletes. Thus 0 ≤ d ≤ n+m. If d is smaller than n+m, can we get a better algorithm? Yes!





The alignment should be inside the 2d+1 band
 ⇒ No need to fill-in the lower and upper triangle
 ⇒ Time complexity: O(dn)





С

2

4

7

7

G

3

3

5

5

4

Example

Α

G

С

Α

Т





$$v(i, j, d) = \max \begin{cases} v(i-1, j-1, d) + s(S[i], S[j]) \\ v(i-1, j, d-1) - \delta & \text{if } d > 0 \\ v(i, j-1, d-1) - \delta & \text{if } d > 0 \end{cases}$$

Exercise: Write down the base cases, the memoized version, and the non-recursive version.

Global Pairwise Alignment: More Realistic Handling of Indels



• So reformulate as follows:

Let g(k) be the indel weight for an indel of k letters. Typically, $g(k) \leq k \cdot g(1)$. Let U and V be two sequences of length n and m. Then their global pairwise alignment can be extracted from the dynamic programming computation of $S_{n,m}$, where

$$S_{0,0} = 0, \quad S_{0,j} = -g(j), \quad S_{i,0} = -g(i)$$
$$S_{i,j} = \max \left\{ \begin{array}{l} S_{i-1,j-1} + s(u'_i, v'_j) \\ \max_{1 \le k \le j} \{S_{i,j-k} - g(k)\} \\ \max_{1 \le k \le i} \{S_{i-k,j} - g(k)\} \end{array} \right\}$$





- $g(q):N \rightarrow \Re$ is the penalty of a gap of length q
- Note g() is subadditive, i.e, $g(p+q) \le g(p) + g(q)$
- If $g(k) = \alpha + \beta k$, the gap penalty is called affine
 - A penalty (α) for initiating the gap
 - A penalty (β) for the length of the gap



- Global alignment of S[1..n] and T[1..m]:
 - Denote V(i, j) be the score for global alignment between S[1..i] and T[1..j]
 - Base cases:
 - V(0, 0) = 0
 - V(0, j) = g(j)
 - V(i, 0) = g(i)



• Recurrence for i>0 and j>0,

ſ

$$V(i, j) = \max \begin{cases} V(i-1, j-1) + \delta(S[i], T[j]) & \text{Match/mismatch} \\ \max_{0 \le k \le j-1} \{V(i, k) + g(j-k)\} & \text{Insert T[k+1..j]} \\ \max_{0 \le k \le i-1} \{V(k, j) + g(i-k)\} & \text{Delete S[k+1..i]} \end{cases}$$





- We need to fill in all entries in the n×m table
- Each entry can be computed in O(max{n, m}) time
- ⇒ Time complexity = O(nm max{n, m})
- ⇒ Space complexity = O(nm)


Variations of Pairwise Alignment

- Fitting a "short" seq to a "long" seq
- Find "local" alignment



 Indels at beginning and end are not penalized



- Find *i*, *j*, *k*, *l*, so that
 S(A) is maximized,
 - A is alignment of $u_i \dots u_j$ and $v_k \dots v_l$







- Given two long DNAs, both of them contain the same gene or closely related gene
 - Can we identify the gene?
- Local alignment problem: Given two strings S[1..n] and T[1..m], among all substrings of S and T, find substrings A of S and B of T whose global alignment has the highest score

Brute-Force Solution Source: Ken Sung



- Algorithm:
 - For every substring A of S, for every substring B of T, compute the global alignment of A and B
 - Return the pair (A, B) with the highest score
- Time:
 - There are n² choices of A and m² choices of B
 - Global alignment computable in O(nm) time
 - In total, time complexity = $O(n^3m^3)$
- Can we do better?





- X is a suffix of S[1..n] if X=S[k..n] for some k≥1
- X is a prefix of S[1..n] if X=S[1..k] for some k≤n
- E.g.
 - Consider S[1..7] = ACCGATT
 - ACC is a prefix of S, GATT is a suffix of S
 - Empty string is both prefix and suffix of S

Dynamic Programming for Local Alignment Problem



- all suffixes A of S[1..i] and
- all suffixes B of T[1..j]
- Then, score of local alignment is

– max_{i,j} V(i ,j)





Smith-Waterman Algorithm

• Basis:

V(i, 0) = V(0, j) = 0

• Recursion for i>0 and j>0:

$$V(i, j) = \max \begin{cases} 0 & \text{Ignore initial segment} \\ V(i-1, j-1) + s(S[i], T[j]) & \text{Match/mismatch} \\ V(i-1, j) - \delta & \text{Delete} \\ V(i, j-1) - \delta & \text{Insert} \end{cases}$$

- Score for match = 2
- Score for insert, delete, mismatch = -1





	_	С	Т	С	Α	Т	G	С
_	0	0	0	0	0	0	0	0
A	0							
С	0							
Α	0							
Α	0							
Т	0							
С	0							
G	0							

- Score for match = 2
- Score for insert, delete, mismatch = -1





	_	С	Т	С	Α	т	G	С
_	0	0	0	0	0	0	0	0
Α	0	0	0	0	2	1	0	0
С	0	2	1	2	1	1	0	2
Α	0	0	1	1	4	3	2	1
A	0	0	0	0	3	3	2	1
Т	0	0	2	1	2			
С								
G								



CAATCG

 C_AT_G

	_	С	Т	С	Α	Т	G	С
_	0	0	0	0	0	0	0	0
Α	0	0	0	0	2	1	0	0
С	0	2+	- 1	2	_ 1	1 🗕	0	2
Α	0	0	1	1	4 ←	- 3 ←	- 2 +	- 1
Α	0	● ←	0	0	3	3	2	1
Т	0	0	2 ←	_ 1	2	5 ←	- 4 +	- 3
С	0	2	_ 1	4 ←	- 3	4	4	6
G	0	1	1	3	3	3	6 +	- 5

Example (III) Source: Ken Sung





- Need to fill in all entries in the n×m matrix
- Each entries can be computed in O(1) time
- Finally, finding the entry with the max value
- \Rightarrow Time complexity = ??
- \Rightarrow Space complexity = O(nm)

Exercise: What is the time complexity?

Multiple Sequence Alignment





What is a domain

- A domain is a component of a protein that is selfstabilizing and folds independently of the rest of the protein chain
 - Not unique to protein products of one gene; can appear in a variety of proteins
 - Play key role in the biological function of proteins
 - Can be "swapped" by genetic engineering betw one protein and another to make chimeras
- May be composed of one, more than one, or not any structural motifs (often corresponding to active sites)

Discovering Domain and Active Sites

>gi|475902|emb|CAA83657.1| protein-tyrosine-phosphatase alpha MDLWFFVLLLGSGLISVGATNVTTEPPTTVPTSTRIPTKAPTAAPDGGTTPRVSSLNVSSPMTTSAPASE PPTTTATSISPNATTASLNASTPGTSVPTSAPVAISLPPSATPSALLTALPSTEAEMTERNVSATVTTQE TSSASHNGNSDRRDETPIIAVMVALSSLLVIVFIIIVLYMLRFKKYKQAGSHSNSFRLPNGRTDDAEPQS MPLLARSPSTNRKYPPLPVDKLEEEINRRIGDDNKLFREEFNALPACPIQATCEAASKEENKEKNRYVNI LPYDHSRVHLTPVEGVPDSHYINTSFINSYQEKNKFIAAQGPKEETVNDFWRMIWEQNTATIVMVTNLKE RKECKCAQYWPDQGCWTYGNIRVSVEDVTVLVDYTVRKFCIQQVGDVTNKKPQRLVTQFHFTSWPDFGVP FTPIGMLKFLKKVKTCNPQYAGAIVVHCSAGVGRTGTFIVIDAMLDMMHAERKVDVYGFVSRIRAQRCQM VQTDMQYVFIYQALLEHYLYGDTELEVTSLEIHLQKIYNKVPGTSSNGLEEEFKKLTSIKIQNDKMRTGN LPANMKKNRVLQIIPYEFNRVIIPVKRGEENTDYVNASFIDGYRRTPTCQPRPVQHTIEDFWRMIWEWK SCSIVMLTELEERGQEKCAQYWPSDGSVSYGDINVELKKEEECESYTVRDLLVTNTRENKSRQIRQFHFH GWPEVGIPSDGKGMINIIAAVQKQQQQSGNHPMHCHCSAGAGRTGTFCALSTVLERVKAEGILDVFQTVK SLRLQRPHMVQTLEQYEFCYKVVQEYIDAFSDYANFK

• How do we find the domain and associated active sites in the protein above?

Domain/Active Sites as Emerging Patterns

- How to discover active site and/or domain?
- If you are lucky, domain has already been modelled
 - BLAST,
 - HMMPFAM, ...
- If you are unlucky, domain not yet modelled
 - Find homologous seqs
 - Do multiple alignment of homologous seqs
 - Determine conserved positions
 - \Rightarrow Emerging patterns relative to background
 - \Rightarrow Candidate active sites and/or domains



In the course of evolution...



Copyright 2007 © Limsoon Wong

Multiple Alignment: An Example

- Multiple seq alignment maximizes number of positions in agreement across several seqs
- seqs belonging to same "family" usually have more conserved positions in a multiple seq alignment

gi 126467
gi 2499753
gi 462550
gi 2499751
gi 1709906
gi 126471
gi 548626
gi 131570
gi 2144715

FHFTSWPDFGVPFTPIGMLKFLKKVKACNP--QYAGAIVVHCSAGVGRTGTFVVIDAMLD FHFTGWPDHGVPYHATGLLSFIRRVKLSNP--PSAGPIVVHCSAGVGRTGCYIVIDIMLD YHYTQWPDMGVPEYALPVLTFVRRSSAARM--PETGPVIVHCSAGVGRTGTTIVIDSMLQ FHFTSWPDHGVPDTTDLLINFRYLVRDYMKQSPPESPILVHCSAGVGRTGCFIVIDAMLE LHFTSWPDHGVPEHPTPFLAFLRRVKTCNP--PDAGPMVVHCSAGVGRTGCFIVIDAMLE LHFTSWPDFGVPFTPIGMLKFLKKVKTLNP--VHAGPIVVHCSAGVGRTGTFIVIDAMLA FHFTGWPDHGVPYHATGLLSFIRRVKLSNP--PSAGPIVVHCSAGAGRTGCYIVIDIMLD FHFTGWPDHGVPYHATGLLGFVRQVKSKSP--PNAGPLVVHCSAGAGRTGCFIVIDIMLD FHFTSWPDHGVPYHATGLLGFVRQVKSKSP--PNAGPLVVHCSAGVGRTGTFIAIDRLIY

Conserved sites

Multiple Alignment: Naïve Approach



 Let S(A) be the score of a multiple alignment A. The optimal multiple alignment A of sequences U₁, ..., U_r can be extracted from the following dynamic programming computation of S_{m1},...,mr:

$$S_{m_1,\dots,m_r} = \max_{\epsilon_1 \in \{0,1\},\dots,\epsilon_r \in \{0,1\}} \left\{ \begin{array}{c} S_{m_1-\epsilon_1,\dots,m_r-\epsilon_r} + \\ s(\epsilon_1 \cdot u'_{1,m_1},\dots,\epsilon_r \cdot u'_{r,m_r}) \end{array} \right\}$$

where

$$\epsilon_i \cdot a = \left\{ egin{array}{cc} a & ext{if } \epsilon_i = 1 \ - & ext{if } \epsilon_i = 0 \end{array}
ight.$$

• This requires O(2^r) steps

Exercise for the Brave: Propose a practical approximation

Popular Tools for Sequence Comparison: FASTA, BLAST, Pattern Hunter





Scalability of Software



- Increasing number of sequenced genomes: yeast, human, rice, mouse, fly, ...
- S/w must be "linearly" scalable to large datasets

Need Heuristics for Sequence Comparison



⇒ Given current size of sequence databases, use of optimal algorithms is not practical for database search

- Heuristic techniques:
 - BLAST
 - FASTA
 - Pattern Hunter
 - MUMmer, ...
- Speed up:
 - 20 min (optimal alignment)
 - 2 min (FASTA)
 - 20 sec (BLAST)

Exercise: Describe MUMer





Basic Idea: Indexing & Filtering

- Good alignment includes short identical, or similar fragments
- ⇒ Break entire string into substrings, index the substrings
- ⇒ Search for matching short substrings and use as seed for further analysis
- ⇒ Extend to entire string find the most significant local alignment segment



- Similarity matching of words (3 aa's, 11 bases)
 - No need identical words
- If no words are similar, then no alignment
 - Won't find matches for very short sequences
- MSP: Highest scoring pair of segments of identical length. A segment pair is locally maximal if it cannot be improved by extending or shortening the segments
- Find alignments w/ optimal max segment pair (MSP) score
- Gaps not allowed
- Homologous seqs will contain a MSP w/ a high score; others will be filtered out



Step 1

• For the query, find the list of high scoring words of length w



Query Sequence of length L

Maximum of L-w+1 words (typically w = 3 for proteins)

For each word from the query sequence find the list of words that will score at least T when scored using a pair-score matrix (e.g. PAM 250).



Step 2

• Compare word list to db & find exact matches





Step 3

• For each word match, extend alignment in both directions to find alignment that score greater than a threshold s





Spaced Seeds

- 111010010100110111 is an example of a spaced seed model with
 - 11 required matches (weight=11)
 - 7 "don't care" positions

• 11111111111 is the BLAST seed model for comparing DNA seqs



- Seed models w/ different shapes can detect different homologies
 - the 3rd base in a codon "wobbles" so a seed like 110110110... should be more sensitive when matching coding regions
- \Rightarrow Some models detect more homologies
 - More sensitive homology search
 - PatternHunter I
- \Rightarrow Use >1 seed models to hit more homologies
 - Approaching 100% sensitive homology search
 - PatternHunter II

Exercise: Why does the 3rd base wobbles?

PatternHunter I



Ma et al., *Bioinformatics* 18:440-445, 2002

• BLAST's seed usually uses more than one hits to detect one homology

 \Rightarrow Wasteful

- Spaced seeds uses fewer hits to detect one homology
- \Rightarrow Efficient

TTGACCTCACC? ||||||||||? TTGACCTCACC? 1111111111 1111111111

1/4 chances to have 2nd hit next to the 1st hit

CAA?A??A?C??TA?TGG? |||?|???|?||? CAA?A??A?C??TA?TGG? 111010010100110111 111010010100110111

1/4⁶ chances to have 2nd hit next to the 1st hit

PatternHunter I



Ma et al., *Bioinformatics* 18:440-445, 2002

Proposition. The expected number of hits of a weight-*W* length-*M* model within a length-*L* region of similarity p is $(L - M + 1) * p^W$

Proof. For any fixed position, the prob of a hit is p^{W} . There are L - M + 1 candidate positions. The proposition follows.



- For *L* = 1017
 - BLAST seed expects (1017 - 11 + 1) * p^{11} = 1007 * p^{11} hits
 - But ~1/4 of these overlap each other. So likely to have only ~750 * p¹¹ distinct hits
 - Our example spaced seed expects (1017 - 18 + 1) * $p^{11} = 1000 * p^{11}$ hits
 - But only 1/4⁶ of these overlap each other. So likely to have ~1000 * p¹¹ distinct hits





PatternHunter I

Proposition. The expected number of hits of a weight-W length-M model within a length-L region of similarity p is $(L - M + 1) * p^W$ Proof. For any fixed position, the prob of a hit is p^W . There are L - M + 1 positions. The proposition follows.



Sensitivity of PatternHunter I



Copyright 2007 © Limsoon Wong



Speed of PatternHunter I



Mouse Genome Consortium used PatternHunter to compare mouse genome & human genome

PatternHunter did the job in a 20 CPU-days --it would have taken BLAST 20 CPU-years!



How to Increase Sensitivity?

- Ways to increase sensitivity:
 - "Optimal" seed
 - Reduce weight by 1
 - Increase number of spaced seeds by 1
- Intuitively, for DNA seq,
 - Reducing weight by 1 will increase number of matches 4 folds
 - Doubling number of seeds will increase number of matches 2 folds
- Is this really so?



How to Increase Sensitivity?

- Ways to increase sensitivity:
 - "Optimal" seed
 - Reduce weight by 1
 - Increase number of spaced seeds by 1



Proposition. The expected number of hits of a weight-W length-M model within a length-L region of similarity p is $(L - M + 1) * p^{W}$

Proof. For any fixed position, the prob of a hit is p^{W} . There are L – M + 1 positions. The proposition follows.

- For *L* = 1017 & *p* = 50%
 - 1 weight-11 length-18 model expects 1000/2¹¹ hits
 - 2 weight-12 length-18 models expect 2 * 1000/2¹² = 1000/2¹¹ hits
 - ⇒ When comparing regions w/ >50% similarity, using 2 weight-12 spaced seeds together is more sensitive than using 1 weight-11 spaced seed!

Exercise: Proof this claim

PatternHunter II Li et al, *GIW*, 164-175, 2003



- Idea
 - Select a group of spaced seed models
 - For each hit of each model, conduct extension to find a homology
- Selecting optimal multiple seeds is NP-hard

- Algorithm to select multiple spaced seeds
 - Let A be an empty set
 - Let s be the seed such that A ∪ {s} has the highest hit probability
 - $A = A \cup \{s\}$
 - Repeat until |A| = K
- Computing hit probability of multiple seeds is NPhard



Sensitivity of PatternHunter II



- Solid curves: Multiple (1, 2, 4, 8,16) weight-12 spaced seeds
- Dashed curves: Optimal spaced seeds with weight = 11,10, 9, 8
- ⇒ "Doubling the seed number" gains better sensitivity than "decreasing the weight by 1"


Expts on Real Data

- 30k mouse ESTs (25Mb) vs 4k human ESTs (3Mb)
 - downloaded from NCBI genbank
 - "low complexity" regions filtered out
- SSearch (Smith-Waterman method) finds "all" pairs of ESTs with significant local alignments
- Check how many percent of these pairs can be "found" by BLAST and different configurations of PatternHunter II



Farewell to the Supercomputer Age of Sequence Comparison!

 Sequence Length
 Blastn
 PatternHunter

 816k vs 580k
 47 sec
 9 sec

 4639k vs 1830k
 716 sec
 44 sec

 20M vs 18M
 out of memory
 13 min

Computer: PIII 700Mhz Redhat 7.1, 16 main memory



Image credit: Bioinformatics Solutions Inc



Copyright 2007 © Limsoon Wong

More Advanced Sequence Comparison Methods

PHI-BLAST Iterated BLAST PSI-BLAST SAM



PHI-BLAST (Pattern-Hit Initiated BLAST)



- protein sequence and
- pattern of interest that it contains

• Output

 protein sequences containing the pattern and have good alignment surrounding the pattern • Impact

 able to detect statistically significant similarity between homologous proteins that are not recognizably related using traditional onepass methods

PHI-BLAST: How it works



Copyright 2007 © Limsoon Wong



PHI-BLAST: IMPACT

Conserved domain or motif under investigation	Pattern ^a	GenBank (30) accession no.	Top non-trivial relevant hit found by PHI-BLAST		Top non-trivial relevant hit found by BLAST	
		of query	Accession no.	E-value	Accession no.	E-value
A. P-loop ATPase domain in apoptosis regulators and plant stress response proteins	[GA]xxxxGK[ST]	231729	2213598	0.038	2961373	4.7
B. ATPase domain in mismatch repair protein MutL, type II topoisomerases, histidine kinases, and HS90 molecular chaperones	hxhxDxGxG	127552	488200	0.017	2495364	1.8
C. Nucleotidyltransferase domain in archaeal tRNA nucleotidyltransferases	DhDhhh	2826366	2650333	0.061	2650333	8.6
D. Motif VI of superfamily II helicases in archaeal homologs of bacterial DNA primases	QxxGRx[GA]R	2128723	2499099	0.54		

ISS



(Intermediate Sequence Search)

- Two homologous seqs, which have diverged beyond the point where their homology can be recognized by a simple direct comparison, can be related through a third sequence that is suitably intermediate between the two
- High score betw A & C, and betw B & C, imply A & B are related even though their own match score is low



ISS: Search Procedure





ISS: IMPACT

(a)

Alignment by FASTA of the sequences of amicyanin and domain 1 of ascorbate oxidase

60 70 80 90 100 Amicyanin MPHNVHFVAGVLGEAALKGPMMKKEQAYSLTFTEAGTYDYHCTPHPFMRGKVVVE :..: . ::. :: Ascorbate Oxidase ILQRGTPWADGTASISQCAINPGETFFYNFTVDNPGTFFYHGHLGMQRSAGLYGSLI 70 80 90 100 110 120

> No obvious match between Amicyanin and Ascorbate Oxidase



ISS: IMPACT

Alignment by FASTA of the sequences of amicyanin and domain 1 of ascorbate oxidase with the intermediate plastocyanin sequence

Amicyanin		10 DKATIPSESPFA	20 AAAEVADGAIV	30 VDIAKMKYETP	40 ELHVKVGDTVTW-IN
PLASTOCYANIN	SLFAVAA	LCVGSFFLSAAPAS	AQTVAI-KMG	ADNGMLAFEPS	TIEIQAGDTVQW-VN
Ascorbate ox	idase SQIRH	IYKWEVEYMFWAPNC	CNENIVMG	I-NGQFPGP	TIRANAGDSVVVELT
50 6 REAMPHNVHFV . :::: : NKLAPHNV-VV ::::: NKLHTEGV-VI	0 70 AGVLG : EG-QP : : HWHGILQRG-TPW2	80 EAALKGPMMKKE : . : ELSHKDLAFSPO : :: ADGTASISQCAINPO	90 CQAYSLTFT SETFEATFS SETFFYNFTVD	100 EAGTYDYHCT : ::: : : : EPGTYTYYCE- .::: . : . NPGTFFYHGHI	110 PHPFMRGKVVVE :: ::::: -PHRGAGMVGKIVVQ .:.:: ::. GMQRSAGLYGSLIVD
		Convincing l via Plasto	homology ocyanin	Pri th m	reviously only is part was atched

PSI-BLAST



(Position-Specific Iterated BLAST)

- Given a query seq, initial set of homologs is collected from db using GAP-BLAST
- Weighted multiple alignment is made from query seq and homologs scoring better than threshold
- Position-specific score matrix is constructed from this alignment

- Matrix is used to search db for new homologs
- New homologs with good score are used to construct new positionspecific score matrix
- Iterate the search until no new homologs found, or until specified limit is reached



SAM-T98 HMM Method

- Similar to PSI-BLAST
- But use HMM instead of position-specific score matrix

Iterated seq. comparisons vs pairwise seq. comparison



Copyright 2007 © Limsoon Wong

Any Question?





Acknowledgements

- Some slides on popular sequence alignment tools are based on those given to me by Bin Ma and Dong Xu
- Some slides on Needleman-Wunsch and Smith-Waterman are based on those given to me by Ken Sung



References

- J. Park et al. "Sequence comparisons using multiple sequences detect three times as many remote homologs as pairwise methods", *JMB*, 284(4):1201--1210, 1998
- J. Park et al. "Intermediate sequences increase the detection of homology between sequences", *JMB*, 273:349--354, 1997
- Z. Zhang et al. "Protein sequence similarity searches using patterns as seeds", *NAR*, 26(17):3986—3990, 1996
- B. Ma et al. "PatternHunter: Faster and more sensitive homology search", *Bioinformatics*, 18:440–445, 2002
- M. Li et al. "PatternHunter II: Highly sensitive and fast homology search", *GIW*, 164—175, 2003
- D. Brown et al. "Homology Search Methods", *The Practical Bioinformatician*, Chapter 10, pp 217—244, WSPC, 2004





- S.F.Altshcul et al. "Basic local alignment search tool", *JMB*, 215:403--410, 1990
- S.F.Altschul et al. "Gapped BLAST and PSI-BLAST: A new generation of protein database search programs", *NAR*, 25(17):3389--3402, 1997
- S.B.Needleman, C.D.Wunsch. "A general method applicable to the search for similarities in the amino acid sequence of two proteins", *JMB*, 48:444—453, 1970
- T.F.Smith, M.S.Waterman. "Identification of common molecular subsequences", *JMB*, 147:195—197, 1981