For written notes on this lecture, please read chapter 10 of *The Practical Bioinformatician*, ond chapter 2 and 5 of *Algorithms in Bioinformatics*.

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

Wong Limsoon







- Dynamic Programming
- String Comparison
- Sequence Alignment
 - Pairwise Alignment
 - Needleman-Wunsch global alignment algorithm
 - Smith-Waterman local alignment algorithm
 - Scoring function
 - Multiple Alignment
- Popular tools

- FASTA, BLAST, Pattern Hunter

What is Dynamic Programming





Knapsack 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 to maximize the total benefit?



Formulation of a solution Source: http://mat.gsia.cmu.edu/classes/dynamic/node6.html

 Intuitively, 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 *k*' of size *w* – *w_i* to fill ...

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

- Where
 - $-w_i$ and b_i be weight and benefit for item j
 - g(w) is max benefit that can be gained from a wpound knapsack

Example: Direct recursive evaluation



6

National Un

of Singapore



"Memoize" to avoid recomputation



Remove recursion: Dynamic program Canageore



int s[]; s[0] := 0; s[1] := 30; s[2] := 65; s[3] = 95; for i := 4 .. w do s[i] := max_j{b_j + s[i - w_j]}; return s[w];



8

Sequence Alignment



Why we compare sequences



• The structure of a protein defines its function

 In order for a protein to have a specific function, it must satisfy specific structural constraints

- Protein evolves → amino acid seq changes → protein structure changes → breaks those structural constraints → protein loses function
- The more similar two proteins' amino acid sequences are, the more likely they come from the same ancestor → the more likely they have the same structure and function



- 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-2 1 SLGSLTIAEPAMIAECKTREEVFCICRRL?DR?? 34 p28sis 61 LARGKRSLGSLSVAEPAMIAECKTRTEVFEISRRLIDRTN 100



Sequence alignment



sequences

Applications of sequence comparison

Infer protein function

 When two protein look similar, we conjecture they come from the same ancestor and inherit the ancestor's function

• Find evolution distance between two species

- Evolution modifies the DNA of species →
 Similarity of their genome correlates with their evolutionary distance
- Help genome assembly
 - Human genome project reconstructs the whole genome based on overlapping info of a huge amount of short DNA pieces

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 amicyanin and domain 1 of ascorbate oxidase

60 70 80 90 100 Amicyanin MPHNVHFVAGVLGEAALKGPMMKKEOAYSLTFTEAGTYDYHCTPHPFMRGKVVVE :: Ascorbate Oxidase ILORGTPWADGTASISOCAINPGETFFYNFTVDNPGTFFYHGHLGMORSAGLYGSLI 70 80 90 100 110 120 No obvious match between Amicyanin and Ascorbate Oxidase

of Singapore

Sequence alignment: Good example

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

D >gi|13476732|ref|NP_108301.1| unknown protein [Mesorhizobium loti]
gi|14027493|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

NUS National University of Singapore

6

Simple-minded probability & score

Alignment:

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 Copyright 2015 © Wong Limsoon

Global pairwise alignment: **Problem definition**



- 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
- 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

Global pairwise alignment: Dynamic programming solution

- Define an indel-similarity matrix s(.,.); e.g.,
 - 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



Needleman-Wunsch algorithm (I Source: Ken Sung

- Consider two strings S[1..n] and T[1..m]
- Let V(i, j) be score of optimal 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 (Il Source: Ken Sung

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



of Singapore





Source:	Ken	Sung
---------	-----	------

	_	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							





Source: Kei	n Sung
-------------	--------

	_	Α	G	С	Α	Т	G	С	
_	0	1 _	_ - 2 _∢	3_	4_	_ – 5₊	_ - 6 _∢	7	
Α	-1	2							
С	-2	$\int S_{0}$),0 +	-s(z)	(A, A)		0	+	2
AS _{1,1}	n = m	$\operatorname{ax} \left\{ S_{0} \right\}$	D,1 —		1 :	= max	x{-1	_	1 = 2
A	- 4	S_1	,0 -		1		[-1		1
Т	- 5								
С	-6								
С	-7								



	Example (III) Source: Ken Sung													
	_	Α	G	С	Α	Т	G	С						
_	0	1 _	_ - 2 _∢	3_	_ – 4	_ – 5,	6₄	7						
Α	-1	2	_ 1											
С	-2	$\int S_0$),1 +	s(A	,G)		[-1	+ -	-1					
AS _{1,1}	$_2 = ma$	$ax \{ S_0 \}$,2 —		1 =	= max	{-2	_	1 =					
A	- 4	$\bigcup S_1$,1		1		2	_	1					
т	- 5													
С	–6													
С	-7													

Copyright 2015 © Wong Limsoon





	_	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. 2									
С	- 7					_, _ , , , ,	(,~),		,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,		



Example (V) Source: Ken Sung





```
Pseudo codes
                        Source: Ken Sung
Create the table V[0..n, 0..m] and P[1..n, 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,j] := V[i,j-1] - \delta;
       set P[i,j] := (0, -1);
       if V[i,j] < V[i-1,j] - \delta then
              set V[i,j] := V[i - 1,j] - \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)
- \Rightarrow Space complexity = O(nm)

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



Problem on speed

- 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!



 $\mathbf{29}$



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





Example



	_	Α	G	С	Α	Т	G	С
	0	1 👞	2 🗸	3				
Α		2	1	_ 0 _	1			
С	-2	1	1	3	2	_ 1		
Α	-3	- 0	•	2	5 ↓	- 4 🗸	- 3	
Α		-1	→ 1	1	4	4	3	_ 2
Т			-2	0	3	6 ←	- 5 +	- 4
С				0	2	5	5	7
С					1	4	4	7



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

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

NUS National University of Singapore

Problem on space

- Dynamic programming requires O(mn) space
- When we compare two very long sequences, space may be the limiting factor
- Can we solve the string alignment problem in linear space?



Easy, if no need to recover alignment

- When filling row 4, it depends only on row 3
 - No need to keep rows 1 and 2
- I.e., we only need to keep two rows

	_	А	G	С	А	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

 \Rightarrow "Cost only" algo

Recovering alignment in O(n+m) space

- Use cost-only algo to find mid-point of alignment
- Divide the problem into two halves
- Recursively deduce alignments for the two halves





How to find mid-point

 $V(S[1..n], T[1..m]) = \max_{0 \le j \le m} \{V(S[1..\frac{n}{2}], T[1..j]) + V(S[\frac{n}{2} + 1..n], T[j+1..m])\}$

- Do cost-only dynamic programming for 1st half
 I.e., find V(S[1..n/2], T[1..j]) for all j
- Do cost-only dynamic programming for 2nd half
 i.e., find V(S[n/2+1..n], T[j+1..m]) for all j
- Determine j which maximizes the sum above



Example Step 2

Step 1

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

Step 4: Recursive on subproblems

	1											
	_	Α	G	С	А	Т	G	С	_			
_												
А												
С												
А												
А												
Т												
С												
С												
_												

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

Step 3




Complexity analysis

Space

- O(m) working memory for finding mid-point
- Once mid-point is found, can free working memory → In each recursive call, we only need to store the alignment path
- Alignment subpaths are disjoint → total space required is O(n+m)
- **Time?** This one is for you to think about \bigcirc



Global pairwise alignment: More Realistic Handling of Indels

- In Nature, indels of several adjacent letters are not the sum of single indels, but the result of one event
- 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)

N-W algo w/ general gap penalty (

Recurrence for i>0 and j>0,

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

National University of Singapore





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



Variations of pairwise alignment ²

- Fitting a "short" seq to a "long" seq
- U V_____
- Indels at beginning and end are not penalized



Find "local" alignment

- 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

Which other string is both a prefix and suffix of S?



Dynamic programming for local alignment problem

- Define V(i, j) be max score of global alignment of A and B over
 - all suffixes A of S[1..i] and
 - all suffixes B of T[1..j]
- Then, score of local alignment is

 max_{i,i} 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
Α	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
Α	0	0	0	0	3	3	2	1
Т	0	0	2	1	2			
С								
G								





	_	С	Т	С	Α	Т	G	С	
_	0	0	0	0	0	0	0	0	A lo
Α	0	0	0	0	2	1	0	0	is
С	0	2 ←	- 1	2	1	1 🔶	0	2	
Α	0	0	1		4 ←	- 3 ←	- 2 +	1	
Α	0	0	0	0	3	3 🔶	2	_ 1	
Т	0	0	2 +	- 1	2	5 ←	- 4 +	3	W
С	0	2	1	4 ←	- 3	4	4	6 ▲	otl loc
G	0	1	1	3	3	3	6 ←	- 5	

An optimal local alignment is

```
C_AT_G
CAATCG
```

What is the other optimal local alignment?





- 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 = ??
- ⇒ Space complexity = O(nm)

Exercise: What is the time complexity?

Recent photos



53

Limsoon & Temple Smith Ken & Michael Waterman





Scoring Function





Scoring function for DNA

- For DNA, since we only have 4 nucleotides, the score function is simple
 - BLAST matrix
 - Transition-transversion matrix: give mild penalty for replacing purine by purine. Similar for replacing pyrimadine by pyrimadine

	Α	С	G	Т
Α	5	-4	-4	-4
С	-4	5	-4	-4
G	-4	-4	5	-4
Т	-4	-4	-4	5

BLAST Matrix



Transition-Transversion Matrix

Copyright 2015 © Wong Limsoon



Scoring function for Protein

- Commonly, it is devised based on two criteria:
 - Chemical/physical similarity
 - Observed substitution frequencies

Scoring function for protein using physical/chemical properties

- An amino acid is more likely to be substituted by another if they have similar property [Karlin & Ghandour, PNAS, 82:8597, 1985]
- The score matrices can be derived based on hydrophobicity, charge, electronegativity, & size
- E.g., give higher score for substituting nonpolar amino acid to another nonpolar amino acid



Scoring function for protein based on statistical model

- Most often used approaches
- Two popular matrices:
 - Point Accepted Mutation (PAM) matrix
 - BLOSUM
- Both methods define the score as the log-odds ratio between the observed substitution rate and the actual substitution rate

Point Accepted Mutation (PAM)



- A point mutation means substituting one residue by another
 - It is called an accepted point mutation if the mutation does not change the protein's function or is not fatal
- Two sequence S1 and S2 are said to be 1 PAM diverged if a series of accepted point mutations can convert S1 to S2 with an average of 1 accepted point mutation per 100 residues

50



PAM matrix by example (I)

- Ungapped alignment is constructed for high similarity amino acid sequences (usually >85%)
- Below is a simplified gap-free global multiple alignment of some highly similar amino acid seqs

IACGCTAFK
 IGCGCTAFK
 LACGCTAFK
 IGCGCTGFK
 IGCGCTLFK
 LASGCTAFK
 LACACTAFK



PAM matrix by example (II)

• Build the phylogenetic tree for the sequences





PAM-1 matrix

$$\delta(a,b) = \log_{10} \frac{O_{a,b}}{E_{a,b}}$$

- O_{a,b} and E_{a,b} are observed and expected freq
 - $O_{a,a} = 99/100$, as PAM-1 assumes 1 mutation per 100 residues
 - For $a \neq b$, $O_{a,b} = F_{a,b} / (100 \Sigma_x \Sigma_y F_{x,y})$ where $F_{a,b}$ is freq of substituting a by b or b by a
 - $E_{a,b} = f_a * f_b$ where f_x is # of x divided by total residues
- E.g., $F_{A,G} = 3$, $F_{A,L}=1$, $f_A = f_G = 10/63$, then $O_{A,G} = 3/(100*2*6) = 0.0025$, $E_{A,G} = (10/63)(10/63) = 0.0252$, $\delta(A,G) = \log (0.0025 / 0.0252) = \log (0.09925) = -1.0034$

PAM-n matrix



63

- Let $M_{a,b} = O_{a,b} / f_a$ be prob that a is mutated to b
- Mⁿ(a,b) is prob that a is mutated to b after n mutations
- PAM-n matrix is created by extrapolating PAM-1
- PAM-n matrix is computed as follows.
 - At time t, suppose the residue is a
 - At time t+1, prob that it becomes j is M(a,b)
 - At time t+2, prob that it becomes j is M²(a,b)
 - ...
 - At time t+n, prob that it becomes j is Mⁿ(a,b)
- \Rightarrow (a,b) entry of PAM-n matrix is log(f_a Mⁿ(a,b)/f_a f_b) = log(Mⁿ(a,b)/f_b)



- PAM did not work well for aligning evolutionarily divergent sequences since the matrix is generated by extrapolation
- Henikoff and Henikoff (1992) proposed BLOSUM
- Unlike PAM, BLOSUM matrix is constructed directly from the observed alignment (instead of extrapolation)



Generating conserved blocks

- In BLOSUM, the input is the set of multiple alignments for nonredundant groups of protein families
- Based on PROTOMAT, blocks of nongapped local aligments are derived
- Each block represents a conserved region of a protein family



Extract frequencies from blocks

- From all blocks, we count the frequency f_a for each amino acid residue a.
- For any two amino acid residues a and b, we count the frequency p_{ab} of aligned pair of a and b.
- For example,
 - ACGCTAFKI GCGCTAFKI ACGCTAFKL GCGCTGFKI GCGCTLFKI ASGCTAFKL ACACTAFKL
- There are 7*9=63 residues, including 9's A and 10's G. Hence, $F_A = 9/63$, $F_G = 10/63$.
- There are $9 * \binom{7}{2} = 189$ aligned residue pairs, including 23 (A,G) pairs. Hence, $p_{AG} = 23 / 189$.



BLOSUM Scoring function

- For each pair of aligned residues a and b, the alignment score $\delta(a,b) = 1/\lambda \ln p_{ab}/(p_a p_b)$
 - p_{ab} is prob that a and b are observed to align together
 - p_a and p_b are freq of residues a and b
 - $-\lambda$ is a normalization constant
- Example: p_L =0.099, p_A =0.074, p_{AL} = 0.0044. With λ =0.347, δ (A,L) = -1.47



What is BLOSUM 62?

- To reduce multiple contributions to amino acid pair freq from the most closely related members of a family, similar seqs are merged within block
- BLOSUM p matrix is created by merging seqs with ≥p% similarity

• Example

- AVAAA, AVAAA, AVAAA, AVLAA, VVAAL
- First 4 seqs have ≥80% similarity. Similarity of last seq with the other 4 sequences is <62%
- For BLOSUM 62, we group first 4 seqs and get $AV[A_{0.75}L_{0.25}]AA$, VVAAL. Then $p_{AV} = 1/5$, $p_{AL} = (0.25 + 1)/5$.



BLOSUM vs PAM

- **BLOSUM 80** ≈ **PAM 1**
- **BLOSUM 62** ≈ **PAM 120**
- **BLOSUM 45** ≈ **PAM 250**
- BLOSUM 62 is the default matrix for BLAST 2.0

Ν Ρ Q R S -2 -1 -1 -1 1 0 D -2 -3 0 -4 -3 -2 0 -1 2 0 0 -2 0 -3 -3 -1 1 3 0 0 -3 -4 -2 -3 -1 -2 -3 -4 -3 0 -2 -2 -2 0 -2 -3 -2 -3 -2 -4 -2 6 -1 -3 -1 0 -1 -2 8 -3 -1 -3 -2 1 -2 0 0 -1 -2 -3 -2 2 -1 -1 -3 -3 0 -4 -3 4 -3 2 1 -3 -3 -3 -3 -2 -1 3 -3 -1 -1 -3 5 -2 -1 0 -1 1 2 0 -1 -2 -3 -2 -1 -3 -1 -1 -4 -3 0 -3 2 -2 4 2 - 3 -2 -2 -2 -1 1 -2 -1 -4 -3 -2 -2 0 -1 -1 -2 -1 -3 1 -1 2 -1 1 -1 -1 0 5 -2 -3 1 0 -3 0 1 -3 0 -3 -2 6 -2 0 0 1 0 -3 -4 -2 -1 -3 -1 -1 -4 -1 -3 -2 -2 -1 -2 -1 -1 -2 -4 -3 -2 -3 -1 -3 0 2 -3 -2 0 -3 1 -2 0 0 -1 5 1 0 -1 -3 -2 0 -3 -2'0-3 2-2-1'0-2 1 5 -1 -1 -3 -3 -2 ດ໌ 0 -1 -2 0 -2 -1 1 -1 0 -1 1 -1 -2 0 -3 -2 s -1 -1 -2 -2 -2 -1 -1 -1 0 -1 -1 1 0 -2 -2 т 0 -1 5 4 -3 -1 0 -1 -3 -2 -1 -3 -3 3 -2 1 1 -3 -2 -2 -3 -2 0 -3 -2 -4 -3 1 -2 -2 -3 -3 -2 -1 -4 -4 -2 -3 -3 -2 -3 -2 -2 -3 -2 3 -3 2 -1 -2 -1 -1 -2 -3 -1 -2 -2 -2 -1 2

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 TSSASHNGNSDRDETPIIAVMVALSSLLVIVFIIIVLYMLRFKKYKQAGSHSNSFRLPNGRTDDAEPQS MPLLARSPSTNRKYPPLPVDKLEEEINRRIGDDNKLFREEFNALPACPIQATCEAASKEENKEKNRYVNI LPYDHSRVHLTPVEGVPDSHYINTSFINSYQEKNKFIAAQGPKEETVNDFWRMIWEQNTATIVMVTNLKE RKECKCAQYWPDQGCWTYGNIRVSVEDVTVLVDYTVRKFCIQQVGDVTNKKPQRLVTQFHFTSWPDFGVP FTPIGMLKFLKKVKTCNPQYAGAIVVHCSAGVGRTGTFIVIDAMLDMMHAERKVDVYGFVSRIRAQRCQM VQTDMQYVFIYQALLEHYLYGDTELEVTSLEIHLQKIYNKVPGTSSNGLEEEFKKLTSIKIQNDKMRTGN LPANMKKNRVLQIIPYEFNRVIIPVKRGEENTDYVNASFIDGYRRRTPTCQPRPVQHTIEDFWRMIWEWK 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



74

In the course of evolution...





Multiple alignment: 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	FHFTSWPDFGVPFTPIGMLKFLKK	WKACNPQYAGAIVVHCS	GVGRTGTFVVIDAMLD
gi 2499753	FHFTGWPDHGVPYHATGLLSFIRR	VKLSNPPSAGPIVVHCS	AGAGRTGCYIVIDIMLD
gi 462550	YHYTQWPDMGVPEYALPVLTFVRR	SSAARMPETGPVI.VHCS	AGVGRTGTYIVIDSMLQ
gi 2499751	FHFTSWPDHGVPDTTDLLINFRYL	VRDYMKQSPPESPII.VHCS	AGVGRTGTFIAIDRLIY
gi 1709906	FQFTAWPDHGVPEHPTPFLAFLRR	VKTCNPPDAGPM <mark>V</mark> VHCS	AGVGRTGCFIVIDAMLE
gi 126471	LHFTSWPDFGVPFTPIGMLKFLKK	VKTLNPVHAGP IVVHCS	AGVGRTGTFIVIDAMMA
gi 548626	FHFTGWPDHGVPYHATGLLSFIRR	VKLSNPPSAGPIVVHCS	AGAGRTGCYIVIDIMLD
gi 131570	FHFTGWPDHGVPYHATGLLGFVRQ	VKSKSPPNAGPLVVHCS	AGAGRTGCFIVIDIMLD
gi 2144715	FHFTSWPDHGVPDTTDLLINFRYL	VRDYMKQSPPESPILVHCS	AGVGRTGTFIAIDRLIY
	* * * * * * * *	* * * *	* * * * * * * * *

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

$$\epsilon_i \cdot a = \begin{cases} a & \text{if } \epsilon_i = 1\\ - & \text{if } \epsilon_i = 0 \end{cases}$$

• This requires O(2^r) steps

Exercise for the Brave: Propose a practical approximation

Copyright 2015 © Wong Limsoon

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



Scalability



78

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





Database search

- Consider a database D of genomic sequences (or protein sequences)
- Given a query string Q,
 - Look for string S in D which is the closest match to the query string Q
 - Two meanings for closest match:
 - S and Q has a semi-global alignment (forgive the spaces at the two ends of Q)
 - S and Q have a local alignment



Goodness of a search algorithm

Sensitivity

- Ability to detect "true positive"
- Measured as the probability of finding the match given the query and the database sequence has only x% similarity

Specificity

- Ability to reject "false positive"
- A good search algorithm should be both sensitive and specific



Need heuristics for sequence comparison

- Time complexity for optimal alignment is O(n²), where n is seq length
- ⇒ Given current size of seq 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



BLAST in 3 steps Altschul et al, *JMB* 215:403-410, 1990

- 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

BLAST in 3 steps



Altschul et al, JMB 215:403-410, 1990

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).





85

BLAST in 3 steps Altschul et al, *JMB* 215:403-410, 1990

Step 2

• Compare word list to db & find exact matches





86

BLAST in 3 steps Altschul et al, *JMB* 215:403-410, 1990

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

1111111111 is the BLAST seed model for comparing DNA seqs



Observations on spaced seeds

- 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?

Copyright 2015 © Wong Limsoon

PatternHunter I



89

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 111111111

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.



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



Implication

- 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

gapore

seeds

likely to

be more

sensitive

& more

efficient



Copyright 2015 © Wong Limsoon

92



93

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?

NUS National University of Singapore

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 * $1000/2^{12} = 1000/2^{11}$ 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



96

• 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 U {s} has the highest hit probability
 - $-A = A \cup \{s\}$
 - Repeat until |A| = K
- Computing hit probability of multiple seeds is NPhard

But see also Ilie & Ilie, "Multiple spaced seeds for homology search", *Bioinformatics*, 23(22):2969-2977, 2007



Q7

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
- ⇒ "Double the seed number" gains better sensitivity than "decrease 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



Copyright 2015 © Wong Limsoon

99



Farewell to Supercomputer Age of sequence comparison!

Computer: PIII 700Mhz Redhat 7.1, 1G main memory

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



Image credit: Bioinformatics Solutions Inc



Copyright 2015 © Wong Limsoon



About the inventor: Ming Li



• Ming Li

- Canada Research Chair
 Professor of
 Bioinformatics,
 University Professor,
 Univ of Waterloo
- Fellow, Royal Society of Canada. Fellow, ACM. Fellow, IEEE

Concluding Remarks





What have we learned?

- General methodology
 - Dynamic programming
- Dynamic programming applications
 - Pairwise Alignment
 - Needleman-Wunsch global alignment algorithm
 - Smith-Waterman local alignment algorithm
 - Multiple Alignment
- Important tactics
 - Indexing & filtering (BLAST)
 - Spaced seeds (Pattern Hunter)

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, Smith-Waterman, and scoring functions are based on those given to me by Ken Sung





- 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
- M. O. Dayhoff, R. M. Schwartz, B. C. Orcutt. A model of evolutionary change in proteins. In M. O. Dayhoff (ed) Atlas of Protein Sequence and Structure, vol 5, suppl 3, pp. 345-352, 1978
- S. Henikoff, J.Henikoff, Amino acid substitution matrices from protein blocks. PNAS, 89(biochemistry): 10915 10919, 1992



References

- 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
- 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