

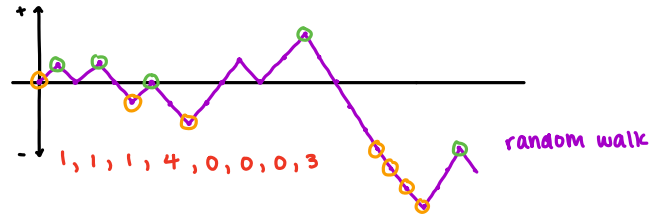
Ch 1: The BLAST Algorithm

1.1 Random Walks

↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
 GGAGACTGTAGACAGCTAATGCTATA
 GAACGCCCTAGCCACGAGCCCTATC
 +1 -1 +1 -1 -1 +1 -1 +1 +1 -1 -1 +1 -1 -1 -1 -1 -1 -1 -1 -1

• ungapped alignment

↳ match +1
 mismatch -1 } for DNA



def: Ladder point ★

• points on the walk lower than any previously reached point

def: Excursion ★

• highest point in the walk from the ladder point before the next ladder point

BLAST

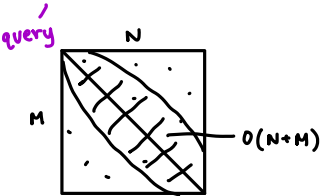
• $O(N+M)$

linear time
empirical
approximation

$J = DB$ - database

$|J| = M$

$|Q| = N$

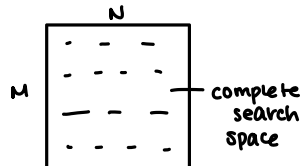


Smith-Waterman

★ both local alignment algorithms

• $O(NM)$

quadratic $N=M$
 $O(N^2)$



Protein Sequences : T Q L A A W C R ...

R H L D D W R R ...

BLOSUM scoring matrix

-1 1 5 -2 1 15 -4 7

• consider an ungapped alignment of two protein sequences both of length N

The Null Hypothesis to be tested is that for each alignment pair of amino acids, the two amino acids were generated by a random process independently such that if amino acid j occurs with probability p_j at any position in the first sequence, and amino acid k occurs at any position in the second sequence with probability p_k , then the probability that they occur together in the alignment is:

$$\text{prob}(j,k) = p_j * p_k$$

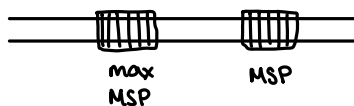
The Alternative Hypothesis

$$\text{prob}(j,k) = q(j,k)$$

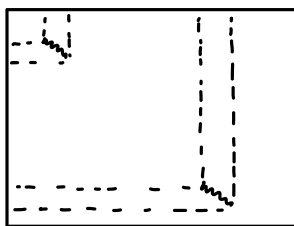
★ $q(j,k)$ function to be determined

$q(j,k)$ related to substitution matrices (BLOSUM, PAM)
(scoring)

Maximum Segment Pair (MSP)



- highest scoring subsequences



BLAST Random Walk

alignment \Rightarrow score: cumulative

Score \Rightarrow random walk

N = number of positions in alignment

$S(j, k)$ = score of aligning amino acid (aa) j with aa k

Scoring Matrix :

two axioms :

AX1: at least one positive score

AX2: The Null Hypothesis score to have negative mean

$$\sum_{j,k} p_j p_k S(j, k) < 0$$

\hookrightarrow when the Null Hypothesis is true, the random walk has a negative drift and will go through a succession of increasingly negative ladder point

Let Y_1, Y_2, \dots be the heights of the excursions of this walk

Let Y_{\max} be the max of these excursions

Y_{\max} = the test statistic of BLAST

It is necessary to find the Null Hypothesis of Y_{\max}

The variables Y_1, Y_2, \dots are identically distributed random variables

The asymptotic distribution of the Y_i is a geometric-like distribution

$$P_r(Y \geq y) \approx c e^{-y^\lambda}$$

* constants c, λ depend on the substitution matrix

p_j, p_k frequencies of aa

m = # ladder points

$$P_r(Y_{\max} \geq y) \approx 1 - e^{-m c e^{-y^\lambda}}$$

$$e^{-m c e^{-y^\lambda}} \leq P_r(Y_{\max} \leq y) \leq e^{-m c e^{-\lambda(y-1)}}$$

P-values for Y_{\max}

$$1 - e^{-m c e^{-y^\lambda}} \leq \text{P-value} \leq 1 - e^{-m c e^{-\lambda(y-1)}}$$