# LECTURE 2. Dynamic Programming & Sequence Alignment
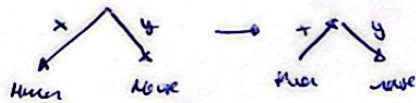
Every organism is thought to come from a common ancestor.

Genome-wide alignments { - highly preserved areas → high functionality
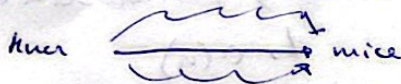
Evolution preserves functional elements

Genomes change over time: mutation, deletion ... S. alignment allows us to find these changes.

- Insertion, deletion, mutation → symmetric operations (i.e., reversibility).



! Exception: CpG dinucleotides aren't symmetric

- Optimality criterion: min number/cost.



- Design algorithm (trade-off)

Example: Formulation 1: Largest Common Substring

S1: A CGT CATCA          offset: 2          A C G T CATCA
S2: TA GT GTCA                              T A G T G TCA

- Algorithm slow
  ↳ make it jump if mismatch → faster

Formulation 2 Largest Common Subsequence
Allow gaps

{ given $X = (x_1, ..., x_m)$, $Z = (z_1, ..., z_k)$ subsequence of $X$ if $i_1 < i_2 < i_3 ... < i_k$



⇒ ~~tess~~ len(LCSS) = 6.
uniform scoring function (= weights)

insertion   deletion   mutation

Formulation 3 — S. Align.
~ Allow gaps (fixed penalty)
  insertion - deletion

- Varying penalties for edit
  transitions (pyrimidine ⟷ pyrimidine, purine ⟷ purine) ↓ cost
  transversion (purine ⟷ pyrimidine) ↑ cost

|   | A | G | T | C |
|---|---|---|---|---|
| A | +1 | -½ | -1 | -1 |
| G | -½ | +1 | -1 | -1 |
| T | -1 | -1 | +1 | -½ |
| C | -1 | -1 | -½ | +1 |

↳ Account for & P(B)

Formulation 4 — Varying gap penalty
  · Linear
  · Affine penalty { 1st nucleotide → p↑
                    : b (don't care about length of mismatch)
  · General → length gap → cost
  · Frame aware { multiple of 3
                  ↳ protein coding regions → less disruptive
  · seek duplications, rearrangements.

How many alignments can be there?

Longest non-boring alignment* (NBA — alignment w gaps always paired up w nucleotides)

$n = len(s1)$
$m = len(s2)$ } $l(less) = n+m$

$2^n$ order problem because @ each position there could be a gap.

↳ we need a polynomial algorithm to find best alignment amongst an exp. n° of alignments → Dynamic Programming

## Intro to DP

Computing Fibonacci seq → $O(2^n)$ { trees nesting inside each other.

def fib(n):
    if n==1 or n==2 : return 1
    return fib(n-1) + fib(n-2)

$T(n) = T(n-1) + T(n-2)\dots$

def fib(n)
    fib_t[1] = 1
    fib_t[2] = 1
    for i in range (3, n+1):
        fib_t[i] = fib_t[i-1] + fib_t[i-2]
    return fib_t[n]

only fills one line of the tree.
↳ $O(n)$

To systematically use a bottom-up approach is successful

Overlapping problems → limited n°
Typically for optimization problems
    ↳ traceback → optimal path
↳ if dependencies between subproblems → no DP

## In practice

Setting up:
    1. Find matrix' parametrization (# directions, variables)
    2. Make sure subproblems are finite
    3. Transversal order (bottom up)
    4. Recursion formula
    5. Remember choices

Start:
    1. Fill results, final op. score
    2. Traceback

# Needleman - Wunch algorithm

Score is additive, smaller to larger:

→ for a given aligned pair $(i,j)$ the best alignment is:

Best of $S1[1...i]$ $S2[1...j]$
+ Best of $S1[i...n]$ $S2[j...m]$ } compute best alignment recursively!

Proof: cut & paste argument

S1 | A | C | G | T | A | T | C | A |

S2 | T | A | G | T | G | T | C | A |

| A | C | G |   | T | C | A | T | A |

| T | A | G | T |   | T | C A |

⇒ Allows for a single recursion (top-left — bottom right) instead of two (middle-to-outside top down)

## Compute score recursively:

### solution #1 : Memoization

o Create a dictionary indexed by aligned seqs when you encounter a new pair of seq

→ if it's in dict: look up the solution

→ if not, compute and add to the dic.

o Ensures no duplication of work!

\# Top down approach.

### Solution #2: Dynamic programming

o Create a big table $(i,j)$

fill it

Explores entire space.
No duplicated work.

\# Bottom up approach

The optimal prefix alt score ⟷ matrix entries
Every optimal algnt. → matrix path
└ best algnt ⟹ best path )

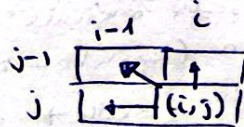| | $S1[1...i]$ | $i$ | $S1[i...n]$ |
|---|---|---|---|
| $S2[0...j]$ | | | |
| $j$ | | | $S$ |
| $S2[j...m]$ | | | |

max. score          diag → match

# DP Approach

- Compute all alignment scores from bottom-up:
  - Define $M[i,j]$ prefix alignment score of $S_1[1\dots i]$ and $S_2[1\dots j]$
  - Fill up table recursively from smaller to bigger alignments.

- Express alignment of $S_1[1\dots i+1]$ and $S_2[1\dots j+1] \rightarrow M[i+1, j+1]$
  - One of three possibilities: (1) extend alignment for $M[i,j]$
    - (2) extend from $M[i-1,j]$ or
    - (3) extend from $M[i, j-1]$ } Only a local computation $O(1)$!

$$M(i,j) = \max \begin{cases} M(i-1,j) - gap \\ M(i-1,j-1) + score \\ M(i,j-1) - gap \end{cases}$$

Only 3 possibilities for extending by one nucleotide: gap in $S_1$, gap in $S_2$, a (mis)match

## Initialization

Create matrix at sizes $(m+1) \times (n+1)$
$$len(s_1) = m$$
$$len(s_2) = n$$

Each $M[i][j]$ represents the optimal alignment score for subsequences $S_1[1\dots i]$ & $S_2[1\dots j]$.

- $M[0][0] = 0$
- $M[0][j] = gap\text{-}penalty * j$    1st first row
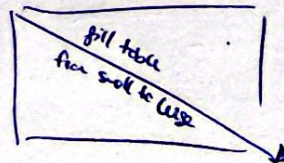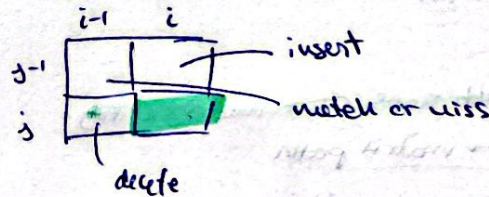- $M[i][0] = gap\text{-}penalty * i$    1st first column

## Fill Matrix

(1) Match or mismatch

match $= M[i-1][j-1] + ($matchscore if $s_1[i-1] == s_2[j-1]$ else mismatch-penalty$)$

(2) delete $= M[i-1][j] + gap$

(3) insert $= M[i][j-1] + gap$

$$M[i][j] = \max(match, insert, delete)$$

↳ contains optimal alignment for subsequences up to

insert / match or miss / delete

fill table from soft to low

# Advanced topics

We can realize that the path tends to be diagonal. Therefore, we can bound the decisions by a diagonal width.

k(N)
↳ width

**Init**

$F(i,0), F(0,j)$ undefined for $i, j > k$

**It.**

for $i = 1 \dots M$

  for $j = \max(1, i-k) \dots \min(N, i+k)$

Can we do better than $O(u^2)$?

→ Probably not.

- - - - -

## Manhattan Tourist Problem

Goal: to find the longest path (i.e. $\max$ ($n°$ attractions or weights))

Vertices:
  ↳ source: northwestern most point
    sink: southeastern most point

source (0,0)

sink (n,m)

- Greedy solution is not optimal!

weight

node

### Approach
- Rather than $(0,0) \to (n,m)$ directly; $(0,0) \to (i,j)$ [arbitrary].
  Best way from source to anywhere?

Find $S_{0,j}$ (for $0 \le j \le m$) is easy   d   $S_{i,0}$ $(0 \le i \le n)$

↳ columns

- Now we compute $S_{1,1}$

    $\overset{(0,1)}{\underset{(1,0) \overset{3}{\to} (1,1)}{\downarrow_0}}$ } $S_{0,1}$ + weight of edge $= \boxed{3}$

$S_{1,0}$ + weight $= 1+3 = \boxed{4}$

- since goal is max. length: $\max \begin{cases} S_{0,1} + edge = 3 \\ S_{1,0} + edge = 4 \end{cases} = 4$

↓

Score for node $(1,1)$

# DAGs

The Manhatten problem can be solved as
a no directed cycles DAG to consider non-perfect grids.

$$G = (V, E)$$

$$E(u, v)$$

- The n° of edges entering a vertex → indegree of v
- The n° of edges leaving a vertex → outdegree of v

$u$ is a predecessor to $v$ if $(u, v) \in E$ —— if it can be reached by travelling backwards.
Therefore $v$ has indegree $k$ if it has $k$ predecessors.

$$S_v = \max_{u \in \text{Predecessors}(v)} (S_u + \text{weight edge}(u, v))$$