

2/18

Announcement : Prof. Kleinberg is holding
extra office hours today from
1:30 pm to 2:30 pm . [Gates 317]

RNA Secondary structure .

RNA molecule: sequence of symbols/bases.

Symbols/bases: $\{A, G, C, U\}$.

A strand of length n :

u_1, u_2, \dots, u_n , each u_i is a base .

Secondary structure with a strand .

\uparrow
 $S \subseteq \{1, \dots, n\} \times \{1, \dots, n\}$

if $(i, j) \in S$ then $i < j$

(i) No sharp turns: if $(i, j) \in S$ then $i < j - 4$.

(ii) if $(i, j) \in S$ then $\{u_i, u_j\}$ is either $\{A, U\}$ or $\{G, C\}$. [valid pairs]

(iii) S to be a matching.

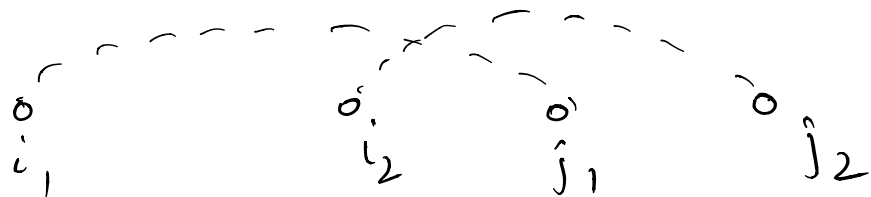
(iv) ^(No crossing) (i_1, j_1) and (i_2, j_2) are in S ,

suppose $i_1 < i_2$, then

either (a) $j_1 < i_2$ or
(b) $j_1 > j_2$.

Not allowed.

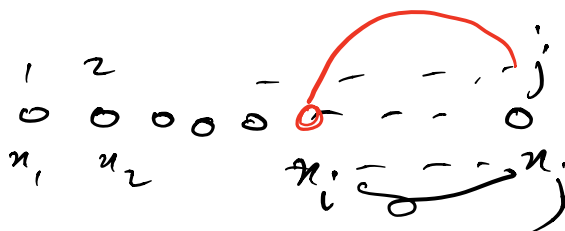
X



Goal: Given RNA strand of length n , find secondary structure of longest cardinality.

Dynamic Programming to solve above problem:

$OPT(j)$: is the optimal secondary structure using the strand $x_1 \dots x_j$.



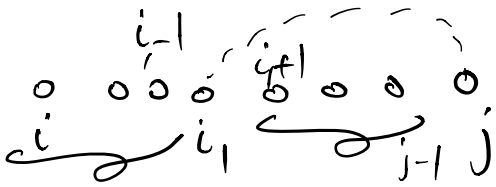
j does not form a secondary bond.
 $OPT(j) = OPT(j-1)$

$OPT(j) = OPT(i-1) + \underbrace{f(i+1, j-1)}_{??} + 1$

leads us to a more general subproblem space:

$OPT(i, j)$: is the optimal secondary structure using the strand $x_i \dots x_j$.

$$OPT(i, j) = \begin{cases} OPT(i, j-1) & \text{(if } j \text{ does not participate)} \\ \text{if } \{x_i, x_j\} \text{ is 'valid'} \\ OPT(i, i_1-1) + OPT(i_1+1, j-1) + 1 \end{cases}$$



$$(*) \quad OPT(i, j) = \max \left\{ \max_{i_1 \in \text{good}(i, j)} \left\{ OPT(i, i_1-1) + OPT(i_1+1, j) + 1 \right\}, OPT(i, j-1) \right\}$$

$$\text{good}(i, j) = \begin{cases} [i, j-5] & \text{if } i < j-4 \\ \emptyset & \text{otherwise} \end{cases}$$

Output $OPT(1, n)$.

Pseudocode:

for $k = 1$ to $n-1$
 for $i = 1$ to $n-k$

 compute $OPT(i, i+k)$ using \otimes
 endfor

endfor
Output $OPT(1, n)$

$O(n^3)$ running time.

Solving recurrences.

easy: $f(0) = c$
 $f(n) = f(n-1) + c'$.

unrolling recursions:

$$\begin{aligned} f(n) &= f(n-1) + c' \\ &= f(n-2) + 2c' \end{aligned}$$

$$f(n) = f(0) + nc' = \Theta(n).$$

$$\rightarrow f(1) = C$$

$$f(n) = f\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + C'$$

unroll $\log n$ times (levels)

$$\begin{aligned} f(n) &= C' \log n + f(1) \\ &= O(\log n) \end{aligned}$$

$$\rightarrow T(n) = 2T(n-1) + 1 ; T(0) = 0.$$

'Guess method'.

say : $T(n) = Cn$ for some constant C .

$$\begin{aligned} T(n) &= 2T(n-1) + 1 \\ &\leq 2C(n-1) + 1 \quad (\text{using induction}) \\ &\quad \neq Cn \end{aligned}$$

say: $T(n) = 2^n - 1$

$$T(n) \leq 2(2^{n-1} - 1) + 1 \\ = 2^n - 1 < 2^n \quad \checkmark$$

$$T(n) = O(2^n)$$

$$\rightarrow T(n) = \sqrt{n} T(\sqrt{n}) + n$$

$$T(n) = n^{3/4} T(n^{1/4}) + 2n$$

$$T(n) = \sqrt{n} \left(n^{1/4} T(n^{1/4}) + \sqrt{n} \right) + n$$

$$T(n) = O(n \log \log n)$$