

CS3000: Algorithms & Data

Jonathan Ullman

Lecture 9:

- Dynamic Programming: Edit Distance, RNA Folding

Oct 5, 2018

Office Hours Sched

Mon
4:30-6:30

Tue
5:00-7:00

Wed
3:00-5:00

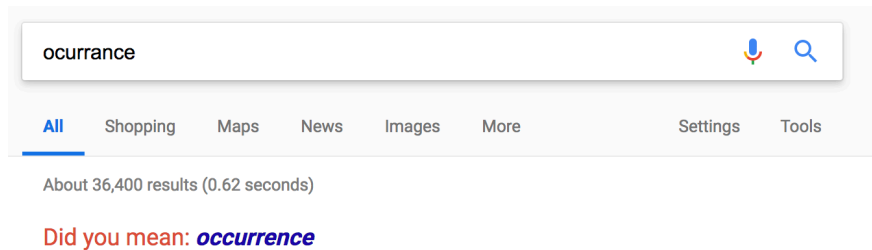
Thu
12:00-2:00
5:00-7:00

} 1 sec
632

Edit Distance Alignments


Distance Between Strings

- Autocorrect works by finding similar strings



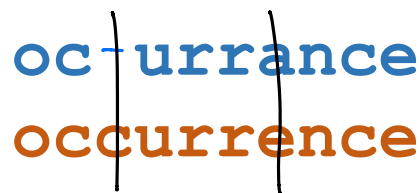
- **ocurrance** and **occurrence** seem similar, but only if we define similarity carefully

ocurrance
occurrence



7 mismatches

oc+urrance
occurrence



2 mismatches

Edit Distance / Alignments

- Given two strings $x \in \Sigma^n$, $y \in \Sigma^m$, the **edit distance** is the number of **insertions**, **deletions**, and **swaps** required to turn x into y .
- Given an **alignment**, the cost is the number of positions where the two strings don't agree

x	o	c		u	r	r	a	n	c	e
y	o	c	c	u	r	r	e	n	c	e

Allow gaps, but not transpositions
cost of the alignment is the # of columns where two symbols don't agree

Ask the Audience

(edit distance)

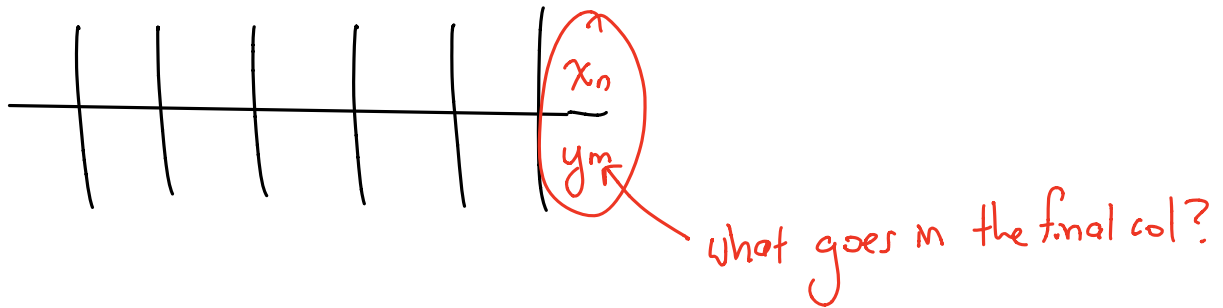
- What is the minimum cost alignment of the strings **smitten** and **sitting** Edit Dist = 3

s	m	i	t	t	e	n	-
s	-	i	t	t	i	n	g
	1				2		3

smitten → sitten → sitt.n → sitting

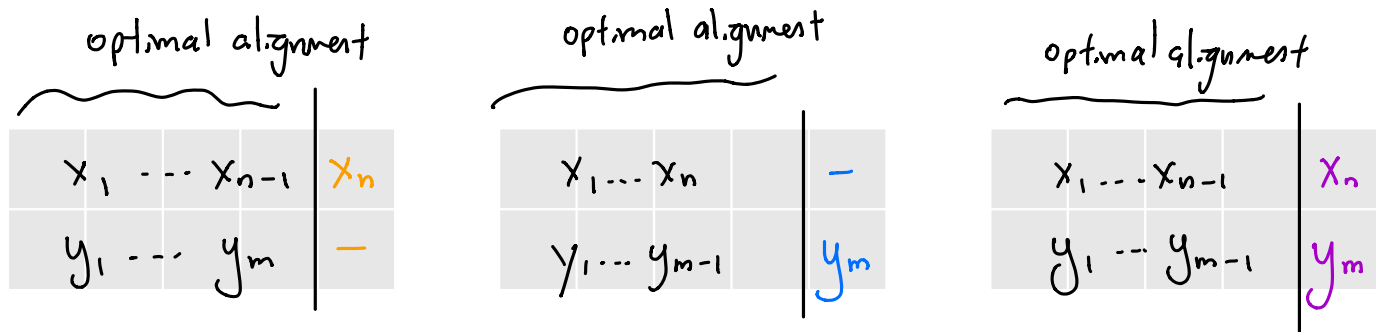
Edit Distance / Alignments

- **Input:** Two strings $x \in \Sigma^n, y \in \Sigma^m$
- **Output:** The minimum cost alignment of x and y
 - **Edit Distance** = cost of the minimum cost alignment



Dynamic Programming

- Consider the **optimal** alignment of x, y
- Three choices for the final column
 - **Case I:** only use x ($x_n, -$)
 - **Case II:** only use y ($-, y_m$)
 - **Case III:** use one symbol from each (x_n, y_m)



To determine which case is best, solve alignment on a smaller x, y .

Dynamic Programming

- Consider the **optimal** alignment of x, y
- **Case I:** only use x ($x_n, -$)
 - deletion + optimal alignment of $x_{1:n-1}, y_{1:m}$
- **Case II:** only use y ($-, y_m$)
 - insertion + optimal alignment of $x_{1:n}, y_{1:m-1}$
- **Case III:** use one symbol from each (x_n, y_m)
 - If $x_n = y_m$: optimal alignment of $x_{1:n-1}, y_{1:m-1}$
 - If $x_n \neq y_m$: mismatch + opt. alignment of $x_{1:n-1}, y_{1:m-1}$

Dynamic Programming

→ $(n+1)(m+1)$ subproblems

- **OPT**(i, j) = cost of opt. alignment of $x_{1:i}$ and $y_{1:j}$
- **Case I:** only use x ($x_i, -$) $1 + \text{OPT}(i-1, j)$
- **Case II:** only use y ($-, y_j$) $1 + \text{OPT}(i, j-1)$
- **Case III:** use one symbol from each (x_i, y_j)

$$\text{OPT}(i, j) = \begin{cases} \text{OPT}(i-1, j-1) + 1 & \text{if } x_i \neq y_j \\ \text{OPT}(i-1, j-1) & \text{if } x_i = y_j \end{cases}$$

$$\text{OPT}(i, j) = \begin{cases} 1 + \min\{\text{OPT}(i-1, j), \text{OPT}(i, j-1), \text{OPT}(i-1, j-1)\} & \text{if } x_i \neq y_j \\ \min\{1 + \text{OPT}(i-1, j), 1 + \text{OPT}(i, j-1), \text{OPT}(i-1, j-1)\} & \text{if } x_i = y_j \end{cases}$$

Dynamic Programming

- **OPT**(i, j) = cost of opt. alignment of $x_{1:i}$ and $y_{1:j}$
- **Case I:** only use x ($x_i, -$)
- **Case II:** only use y ($-, y_j$)
- **Case III:** use one symbol from each (x_i, y_j)

Recurrence:

$$\text{OPT}(i, j) = \begin{cases} 1 + \min\{\text{OPT}(i-1, j), \text{OPT}(i, j-1), \text{OPT}(i-1, j-1)\} \\ \min\{1 + \text{OPT}(i-1, j), 1 + \text{OPT}(i, j-1), \text{OPT}(i-1, j-1)\} \end{cases}$$

Base Cases:

$$\text{OPT}(i, 0) = i, \text{OPT}(0, j) = j$$

Example

x = pert

y = beast

$\begin{array}{|c|} \hline p \\ \hline b \\ \hline \end{array}$

$\begin{array}{|c|c|c|c|c|} \hline opt & p & e & - & r & t \\ \hline & b & e & a & s & t \\ \hline \end{array}$

Edit("p", "bea")

$$OPT(1,1) = 1 + \min\{OPT(0,1), OPT(1,0), OPT(0,0)\}$$

	-	b	e	a	s	t
-	0	1	2	3	4	5
p	1	1	2	3	4	5
e	2	2	1	2	3	4
r	3	3	2	2	3	4
t	4	4	3	3	3	3

$\begin{array}{|c|c|c|} \hline p & \bullet & - \\ \hline b & e & a \\ \hline \end{array}$

Finding the Alignment

- **OPT**(i, j) = cost of opt. alignment of $x_{1:i}$ and $y_{1:j}$
- **Case I:** only use x ($x_i, -$)
- **Case II:** only use y ($-, y_j$)
- **Case III:** use one symbol from each (x_i, y_j)

Edit Distance (“Bottom-Up”)

```
// All inputs are global vars
FindOPT(n,m):
  M[0,j] ← j, M[i,0] ← i

  for (i= 1,...,n):
    for (j = 1,...,m):
      if (xi = yj):
        M[i,j] = min{1+M[i-1,j],1+M[i,j-1],M[i-1,j-1]}
      elseif (xi != yj):
        M[i,j] = 1+min{M[i-1,j],M[i,j-1],M[i-1,j-1]}

  return M[n,m]
```

$[O(nm) \text{ entries}] \times [O(1) \text{ per entry}] = O(nm) \text{ time}$

$O(nm) \text{ space (just for the table)}$

Ask the Audience

- Suppose **inserting/deleting costs $\delta > 0$** and **swapping $a \leftrightarrow b$ costs $c_{a,b} > 0$**
- Write a recurrence for the min-cost alignment

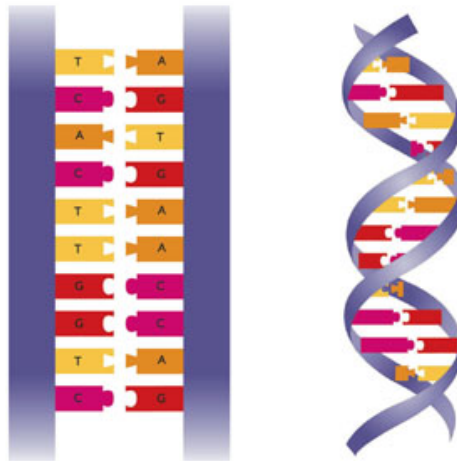
Edit Distance Summary

- Compute the **edit distance**, or **min-cost alignment** between two strings in time/space $O(nm)$
- Dynamic Programming:
 - Decide the final pair of symbols in the alignment
- Space can be prohibitive in practice
 - Compute edit distance in space $O(\min\{n, m\})$
 - Can also find alignment in space $O(n + m)$ using a clever divide-and-conquer algorithm!

RNA Folding

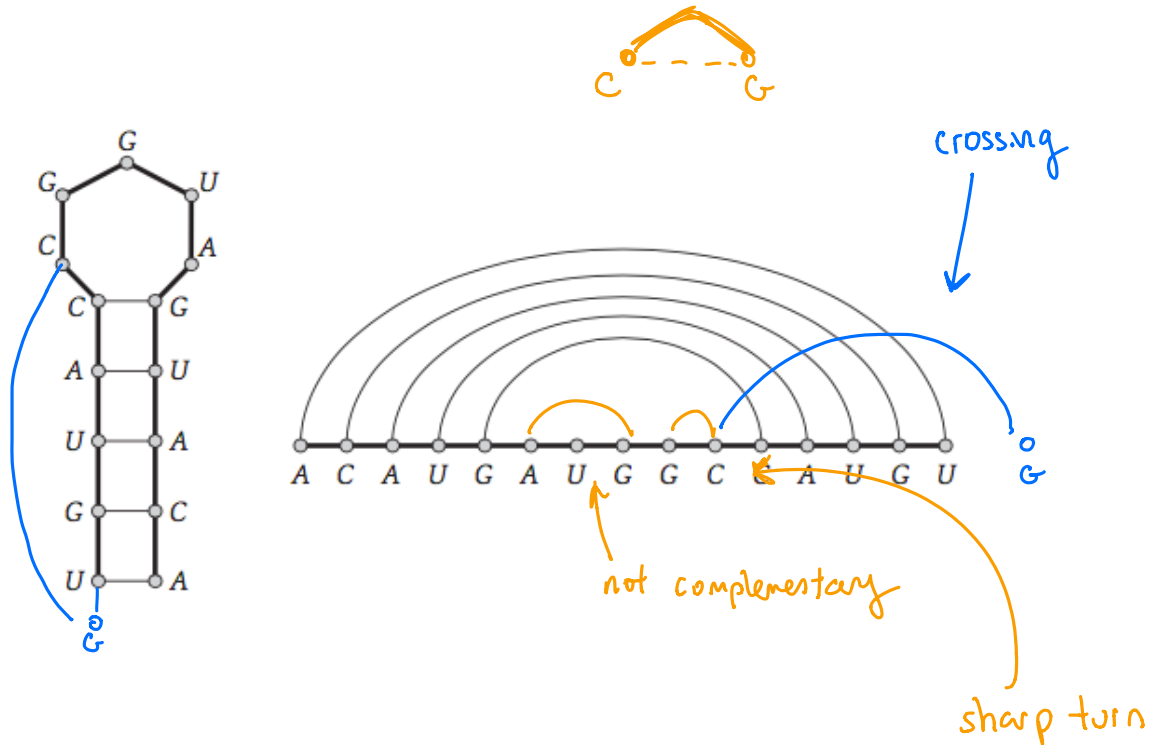
DNA

- DNA is a string of four bases {**A,C,G,T**}
- Two complementary strands of DNA stick together and form a **double helix**
 - **A—T** and **C—G** are complementary pairs



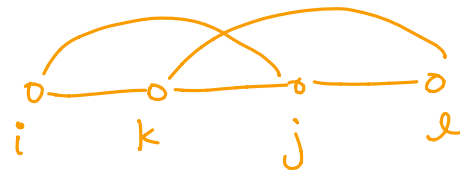
RNA Folding

- RNA strand will try to **minimize energy** (form the most bonds) subject to **constraints**



RNA Folding

- RNA is a string of bases $b_1, \dots, b_n \in \{A, C, G, U\}$
- The structure is given by a set of **bonds** S consisting of pairs (i, j) with $i < j$
 - **(Complements)** Only $A - U$ or $C - G$ can be paired
 - **(Matching)** No base b_i is in two pairs in S
 - **(No Sharp Turns)** If $(i, j) \in S$, then $i < j - 4$
 - **(Non-Crossing)** If $(i, j), (k, \ell) \in S$ then it cannot be the case that $i < k < j < \ell$



RNA Folding

- **Input:** RNA sequence $b_1, \dots, b_n \in \{A, C, G, U\}$
- **Output:** A set of pairs $S \subseteq \{1, \dots, n\} \times \{1, \dots, n\}$
 - **Goal:** maximize the size of S
 - **(Complements)** Only $A - U$ or $C - G$ can be paired
 - **(Matching)** No base b_i is in two pairs in S
 - **(No Sharp Turns)** If $(i, j) \in S$, then $i < j - 4$
 - **(Non-Crossing)** If $(i, j), (k, \ell) \in S$ then it cannot be the case that $i < k < j < \ell$

Dynamic Programming

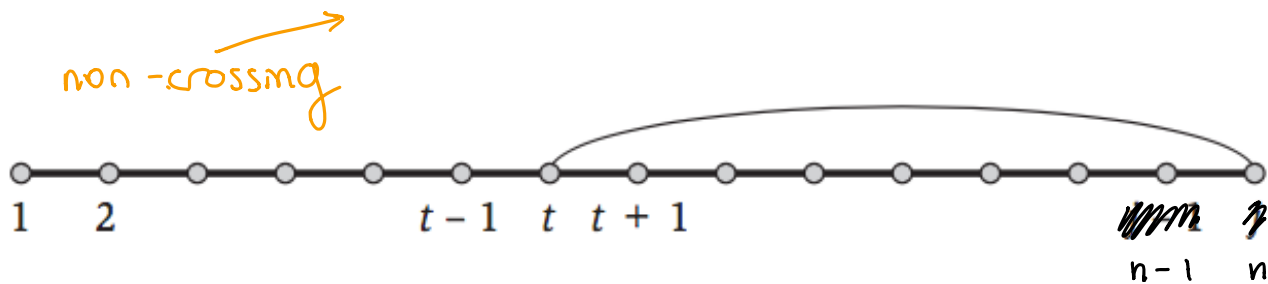
- Let O be the optimal set of pairs for $b_1 \cdots b_n$

- **Case 1:** n pairs with nothing in O

Then O is the opt. set of pairs for b_1, \dots, b_{n-1}

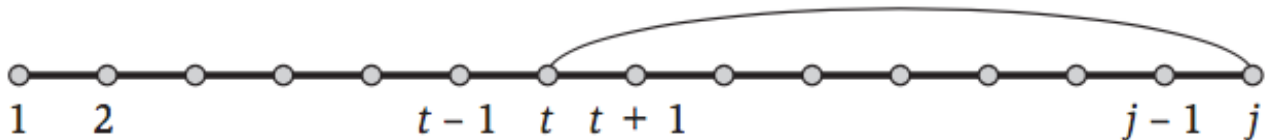
- **Case 2:** n pairs with some $t < n - 4$ in O

Then O is $(t, n) +$ opt. set of pairs for $b_{t+1} \cdots b_{n-1}$
+ opt. set of pairs for $b_1 \cdots b_{t-1}$



Dynamic Programming

- Let $O_{i,j}$ be the optimal set of pairs for $b_i \cdots b_j$
- **Case 1:** j pairs with nothing in $O_{i,j}$
- **Case 2:** j pairs with some $t < j - 4$ in $O_{i,j}$



Dynamic Programming

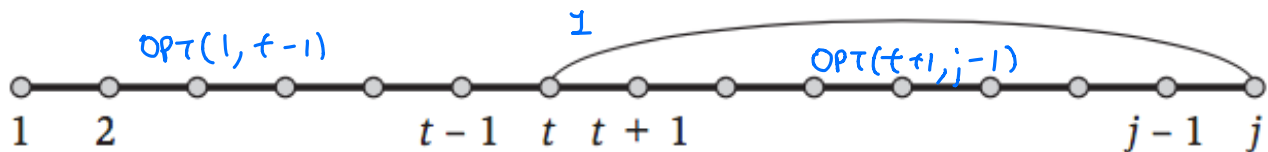
- Let $\text{OPT}(i, j)$ be the opt. **number** of pairs for $b_i \cdots b_j$
- **Case 1:** j pairs with nothing in $O_{i,j}$

$$\text{OPT}(i, j) = \text{OPT}(i, j-1)$$

- **Case 2t:** j pairs with $t < j - 4$ in $O_{i,j}$

- $\text{OPT}(i, j) = 1 + \text{OPT}(i, t-1) + \text{OPT}(t+1, j-1)$

- Consider any t s.t. $t < j-4$ and b_t, b_j are complements



Dynamic Programming

- Let $\text{OPT}(i, j)$ be the opt. **number** of pairs for $b_i \cdots b_j$
- **Case 1:** j pairs with nothing in $O_{i,j}$
- **Case 2t:** j pairs with $t < j - 4$ in $O_{i,j}$

Recurrence:

$$\text{OPT}(i, j)$$

$$= \max\{\text{OPT}(i, j - 1), \max\{\overset{+1}{\text{OPT}(i, t - 1)} + \text{OPT}(t + 1, j - 1)\}\}$$

Maximum over all t such that

- $i \leq t < j - 4$
- b_t, b_j are compatible bases

Base Cases:

$$\text{OPT}(i, j) = 0 \text{ if } i \geq j - 4$$

Filling the Table

Sequence: *ACCGGUAGU*

Recurrence:

$$\text{OPT}(i, j) = \max \left\{ \text{OPT}(i, j - 1), \max_{\text{possible } t} \{ \text{OPT}(i, t - 1) + \text{OPT}(t + 1, j - 1) \} \right\}$$

n-5

	6	7	8	<i>j = 9</i>
4	0	0	0	
3	0	0		
2	0			
<i>i = 1</i>				

n-5

$\Theta(n^2)$ entries
 $\times O(n)$ per entry

RNA Folding Summary

- Compute the **optimal RNA folding** in time $O(n^3)$ and space $O(n^2)$
- **Dynamic Programming:**
 - Decide on an optimal pair $b_t - b_n$
 - Remaining RNA is two non-overlapping pieces
 - **Adding variables:** one subproblem for each interval
- **Non-crossing** and **matching** are critical
 - Think about how the dynamic programming algorithm changes if we remove each of the conditions