

CS3000: Algorithms & Data

Paul Hand

Lecture 12:

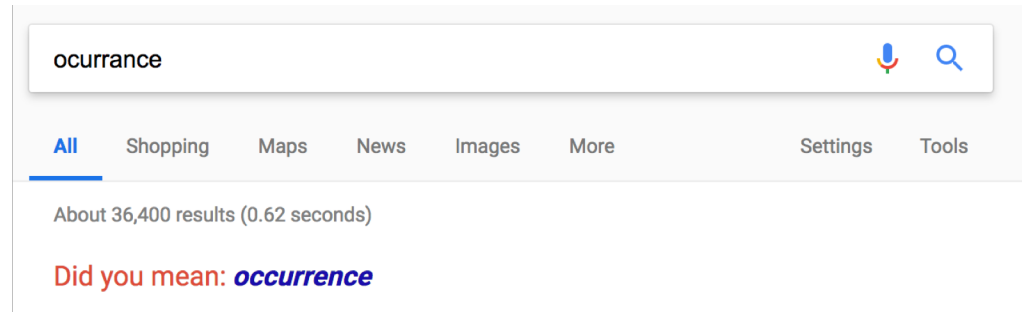
- Dynamic Programming – Sequence Alignment
- Introduction to Graphs

Feb 25, 2019

Sequence Alignments and Edit Distance

Distance Between Strings

- Autocorrect works by finding similar strings



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

ocurrance
occurrence

oc urrance
occurrence

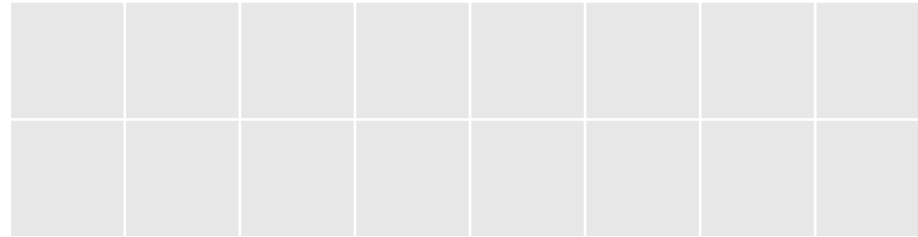
Edit Distance / Alignments

- 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

o	c		u	r	r	a	n	c	e
o	c	c	u	r	r	e	n	c	e

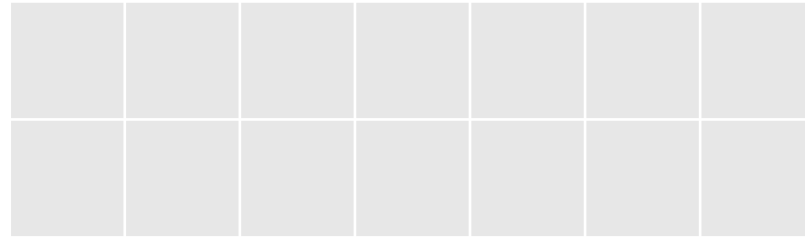
Ask the Audience

- What is the minimum cost alignment of the strings **smitten** and **sitting**



Activity

- Find two strings where two different alignments (insertions, deletions, replacements) realize the edit distance between them.

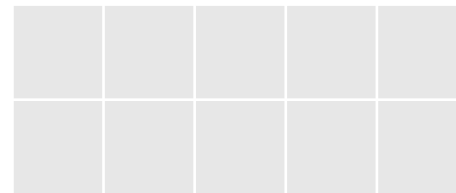


Edit Distance / Alignments

- **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)



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

- **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)

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:

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

Base Cases:

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

Example

x = pert

y = beast

	-	b	e	a	s	t
-						
p						
e						
r						
t						

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]
```

Activity

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

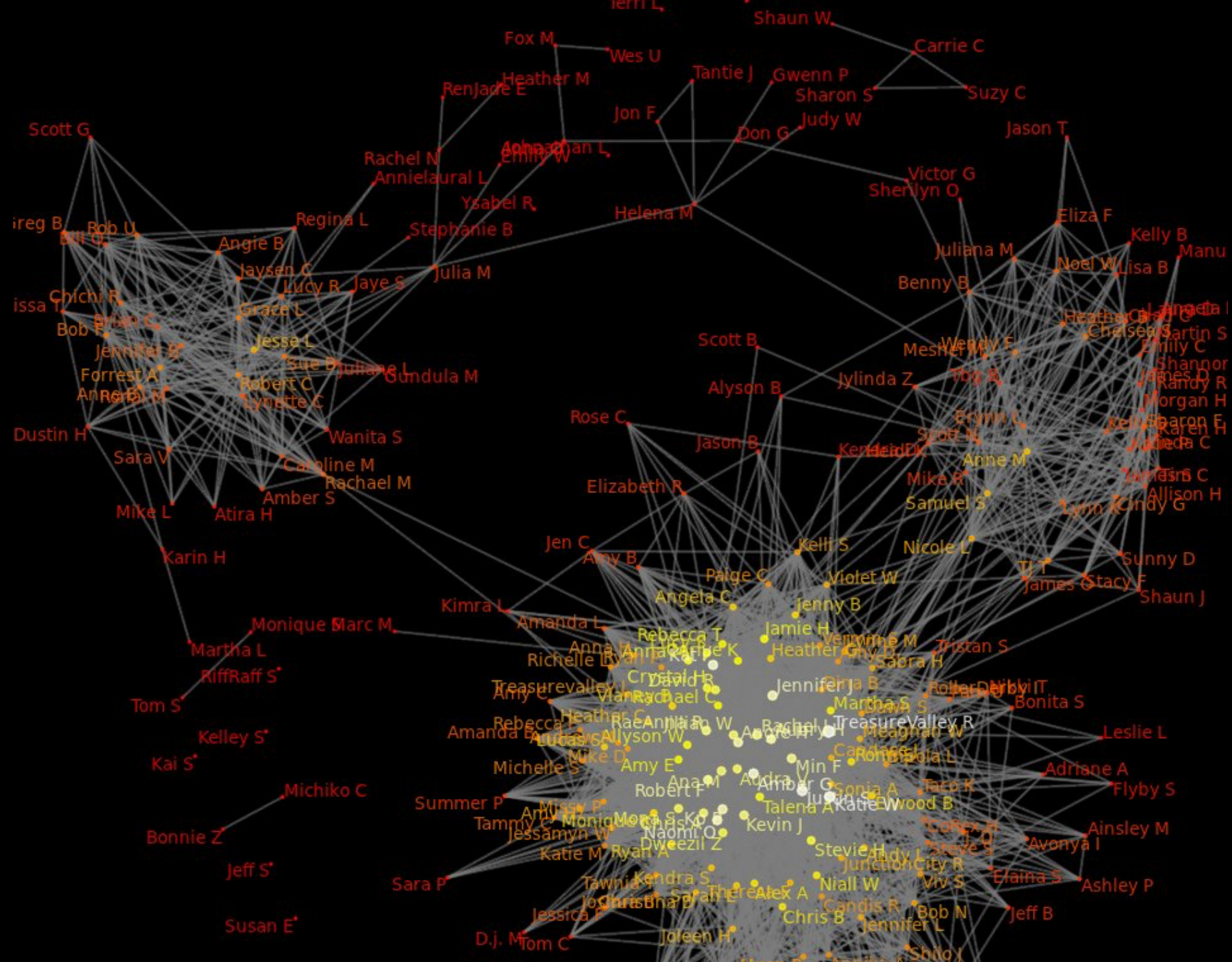
$$\text{OPT}(i, j) = \begin{cases} \min\{1 + \text{OPT}(i - 1, j), 1 + \text{OPT}(i, j - 1), \text{OPT}(i - 1, j - 1)\} & x_i = y_j \\ \min\{1 + \text{OPT}(i - 1, j), 1 + \text{OPT}(i, j - 1), 1 + \text{OPT}(i - 1, j - 1)\} & x_i \neq y_j \end{cases}$$

Discussion

- Dynamic Programming is a time-space tradeoff. Comment on the tradeoff in the case of edit distance.

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

Graphs



Graphs Are Everywhere

- Transportation networks
- Communication networks
- WWW
- Biological networks
- Citation networks
- Social networks
- ...

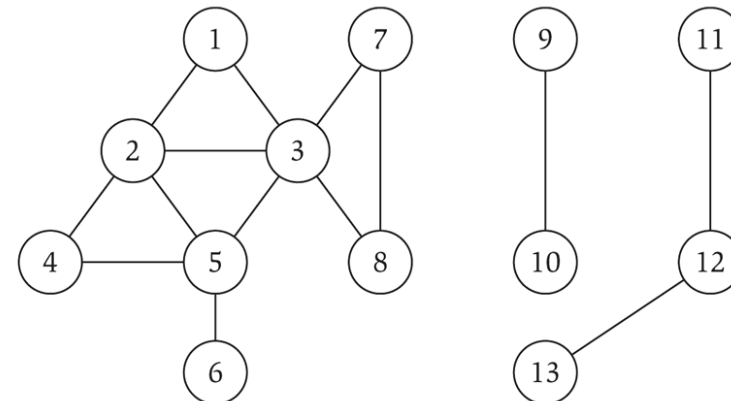
What's Next

- **Graph Algorithms:**
 - **Graphs:** Key Definitions, Properties, Representations
 - **Exploring Graphs:** Breadth/Depth First Search
 - Applications: Connectivity, Bipartiteness, Topological Sorting
 - **Shortest Paths:**
 - Dijkstra
 - Bellman-Ford (Dynamic Programming)
 - **Minimum Spanning Trees:**
 - Borůvka, Prim, Kruskal
 - **Network Flow:**
 - Algorithms
 - Reductions to Network Flow

Graphs: Key Definitions

- **Definition:** A directed graph $G = (V, E)$
 - V is the set of nodes/vertices
 - $E \subseteq V \times V$ is the set of edges
 - An edge is an ordered $e = (u, v)$ “from u to v ”
- **Definition:** An undirected graph $G = (V, E)$
 - Edges are unordered $e = (u, v)$ “between u and v ”

- **Simple Graph:**
 - No duplicate edges
 - No self-loops $e = (u, u)$



Activity

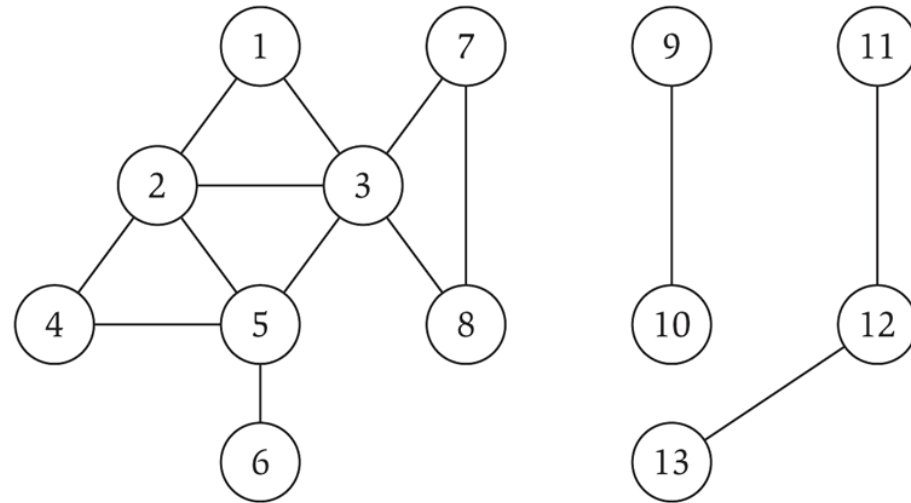
- How many edges can there be in a **simple** **directed/undirected** graph?

Paths/Connectivity

- A **path** is a sequence of consecutive edges in E
 - $P = \{(u, w_1), (w_1, w_2), (w_2, w_3), \dots, (w_{k-1}, v)\}$
 - $P = u - w_1 - w_2 - w_3 - \dots - w_{k-1} - v$
 - The **length** of the path is the # of edges
- An **undirected** graph is **connected** if for every two vertices $u, v \in V$, there is a path from u to v
- A **directed** graph is **strongly connected** if for every two vertices $u, v \in V$, there are paths from u to v and from v to u

Cycles

- A **cycle** is a path $v_1 - v_2 - \dots - v_k - v_1$ where $k \geq 3$ and v_1, \dots, v_k are distinct



Activity: how many cycles are there in this graph?

Activity

- Suppose an undirected graph G is connected
 - True/False? G has at least $n - 1$ edges

Activity

- Suppose an undirected graph G has $n - 1$ edges
 - True/False? G is connected