

Approximate String Matching

* String defn.

$$\bar{X} = x_0 x_1 \dots x_{m-1}, \quad \bar{Y} = y_0 y_1 \dots y_{n-1}$$

* Edit costs

$$\left. \begin{array}{l} \text{match} \\ \text{substitution} \end{array} \right\} \begin{array}{l} x_i = y_j \\ x_i \neq y_j \end{array} \quad d(x_i, y_j)$$

$$\text{deletion} \quad x_i \sim \epsilon \quad d(x_i, \epsilon)$$

$$\text{insertion} \quad \epsilon \sim y_j \quad d(\epsilon, y_j)$$

Symbols $x_i, y_j \in \text{alphabet}$, Symbol ϵ unobservable
(bookkeeping only)

* String alignment, cost

$$\begin{array}{ccccccc} x_0 & - & x_1 & x_2 & \dots & x_{m-1} & \\ y_0 & y_1 & - & y_2 & \dots & y_{n-1} & \\ \uparrow & \uparrow & \uparrow & \uparrow & & \uparrow & \\ \text{match/sub} & & & \text{deletion} & & \text{match/sub} & \\ & & & \text{insertion} & & & \end{array}$$
$$d(\bar{x}, \bar{y}) = \sum_{\text{edit}} d(x_i, y_j)$$

* Optimal string alignment

$$d(\bar{x}, \bar{y}) = \min \sum d(x_i, y_j)$$

all possible
edit seqs.

* Solution

$$d(\epsilon x_0 x_1 \dots x_{m-1}, \epsilon y_0 y_1 \dots y_{n-1}) = \min \{$$

$$\begin{aligned} \text{DEL} \quad & d(\epsilon x_0 x_1 \dots x_{m-2}, \epsilon y_0 y_1 \dots y_{n-1}) + d(x_{m-1}, \epsilon) , \\ \text{MATCH/SUB} \quad & d(\epsilon x_0 x_1 \dots x_{m-2}, \epsilon y_0 y_1 \dots y_{n-2}) + d(x_{m-1}, y_{n-1}) , \\ \text{INS} \quad & d(\epsilon x_0 x_1 \dots x_{m-1}, \epsilon y_0 y_1 \dots y_{n-2}) + d(\epsilon, y_{n-1}) \} \end{aligned}$$

1. Apply recursively \Rightarrow lot of overlapped computation
2. Use memoization to store results (matrix cache)
3. Switch to iteration (forward computation)
4. Reduce mem use if necessary (vector cache)

Cost matrix: $D[m+1][n+1]$

$$D(0,0) = 0 \quad D(i,0) = D(i-1,0) + d(x_i, \epsilon) \quad \text{DEL}$$

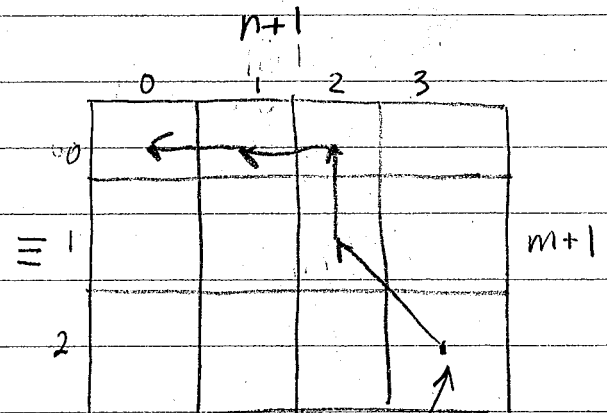
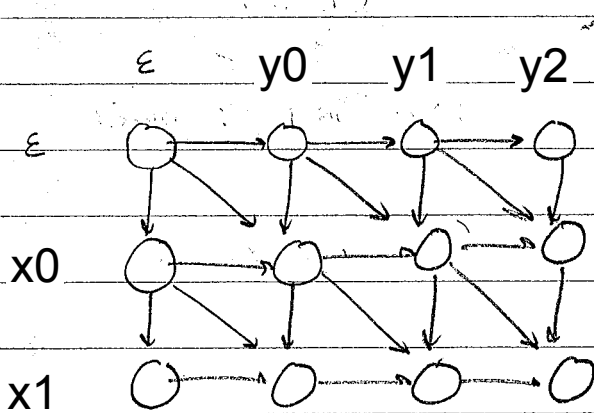
$$D(0,j) = D(0,j-1) + d(\epsilon, y_j) \quad \text{INS}$$

$$D(i,j) = \min \begin{cases} D(i-1,j) + d(x_i, \epsilon) & \text{DEL} \\ D(i-1,j-1) + d(x_i, y_j) & \text{MATCH/SUB} \\ D(i,j-1) + d(\epsilon, y_j) & \text{INS} \end{cases}$$

link matrix: $L[m+1][n+1]$

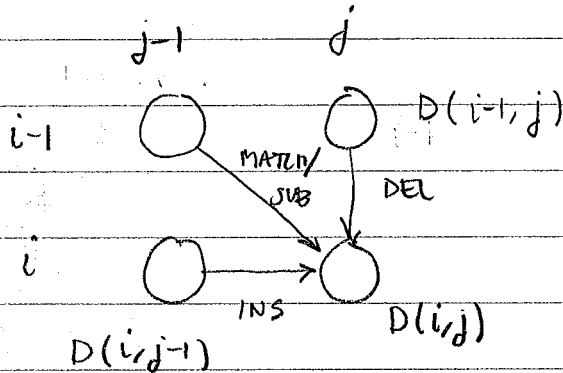
$$L(i,j) = \begin{matrix} \text{MATCH/SUB} & \text{DEL} & \text{diag} & \text{INS} \\ & \text{VERT} & & \text{HORZ} \end{matrix}$$

Graph view of computation



$$D(m,n) = d(\bar{x}, \bar{y})$$

record optimal move(s) in link matrix



Levenshtein edit cost

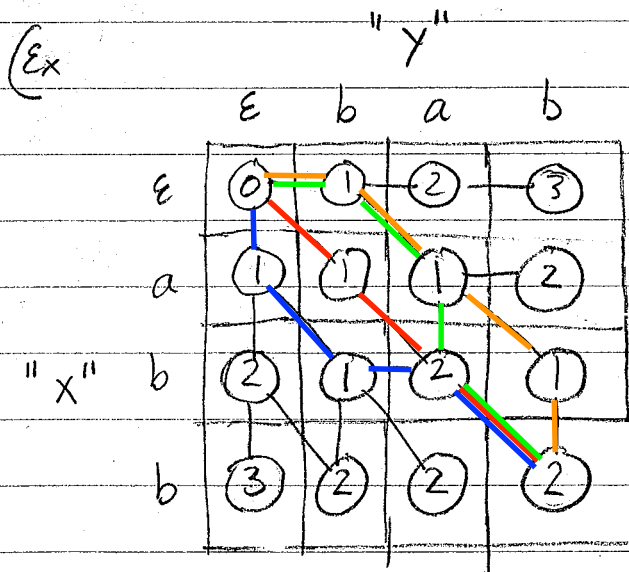
$$d(x_i, y_j) = \begin{cases} 0 & \text{MATCH} \\ 1 & \text{SUB} \end{cases}$$

$$d(x_i, \epsilon) = 1 \quad \text{DEL}$$

$$d(\epsilon, y_j) = 1 \quad \text{INS}$$

(counts num. edits)

choose in order of
DIAG, VERT, HORIZ



- ① a b b ② - a b b ③ a b - b ④ - a b b
 b a b b a - b - b a b b a b -

obtain alignment using backtracking and stack.

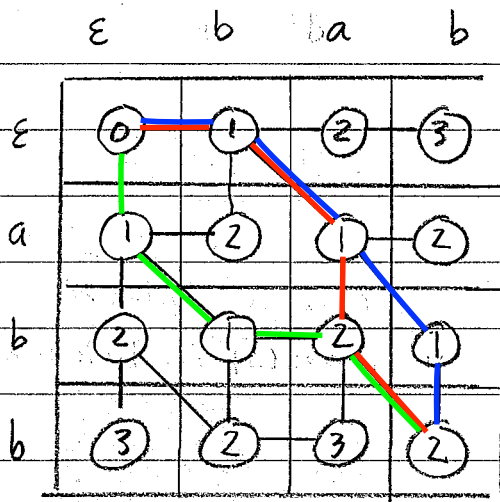
Shortest path thru edit graph

* Longest Common Subsequence (LCS)

$$d(x_i, y_j) = \begin{cases} 0 & \text{MATCH} \\ \infty & \text{SUB} \end{cases} \Rightarrow \text{(SUB not allowed) MATCH}$$

$$d(x_i, \epsilon) = 1 \quad \text{DEL}$$

$$d(\epsilon, x_j) = 1 \quad \text{INS}$$



(2) (1) (3)
 a b - b - a b b - a b b
 - b a b b a - b b a b -

choose in order of
DIAG, VERT, HORIZ

NOTE: \Downarrow

$$D(m, n) = m + n - 2 |LCS|$$

$$|LCS| = (m + n - D(m, n)) / 2$$

(Ex 2 = (3 + 3 - 2) / 2

conv. to the two matches symbol

* Lab 8 : Diff