# String Alignment

## (EDIT DISTANCE)

teh    the

Given: two strings $x, y$ over same alphabet $\Sigma$

operations    $op_1, \ldots, op_k$
  costs    $cost_1, \ldots, cost_k$

Goal: Find the min-cost sequence of operations transforming $x$ into $y$.

For today: our operations will be

substitution          insertion       deletion
(single-letter)

e.g.   "go" → "so"

cost 1                cost 1            cost 1

e.g.    $x =$ THEIR        $y =$ THERE

```
THEIR
|| |ss        cost 2
THERE
```



```
THEIR ⊖ ← not part of alphabet
||| d| i      cost 2
THE R E
```

optimal substructure property ✓
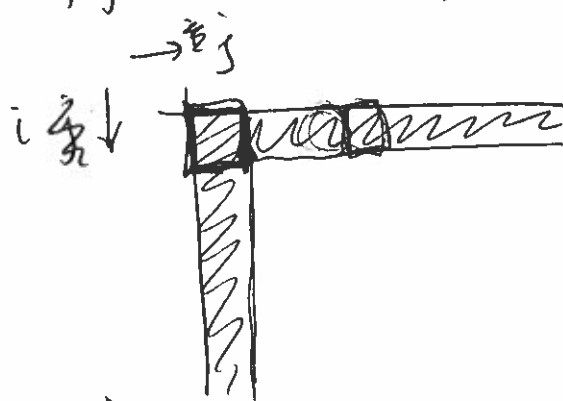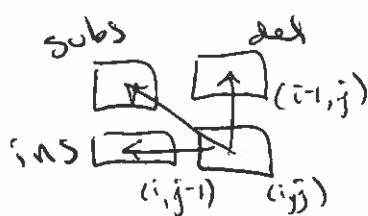
overlapping sub-problems property

Proposed subproblems $(i,j)$ align $x[1..i]$ w/ $y[1..j]$.

$cost(i,j) = $ min cost alignment of $x[1..i]$ w/ $y[1..j]$.

$$= \min \begin{cases} cost(i-1, j-1) + cost(subs) \\ cost(i, j-1) + cost(ins) \\ cost(i-1, j) + cost(del) \end{cases}$$



$cost(0,0) = 0$    (ok base case)

(Note: $cost(0,j) = $ cost to align "" w/ $y[1..j]$.

Also works to say $cost(0,j) = j$    $cost(i,0) = i$.

align (x,y):

~~table~~ n = len(x)
m = len(y)

table (n+1)×(m+1) ~~size~~ array

aux " " "

// Base cases
for i=0 to n
    table $(i,0) = i$
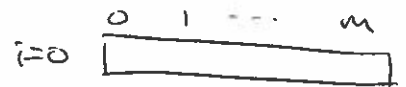    aux $(i,0) = (i-1, 0)$    (if $i \geq 1$)

for j=0 to m
    table $(0, j) = j$
    aux $(0, j) = (0, j-1)$    (if $j \geq 1$)

for i=1 to n
    for j=1 to m
        ~~table~~ poss1 = table $(i-1, j-1) + 1$    // subs
        poss2 = table $(i, j-1) + 1$    // ins
        poss3 = table $(i-1, j) + 1$    // del
        table $(i,j)$ = min (poss1, poss2, poss3)
        aux $(i,j)$ = if poss1 was min then $(i-1, j-1)$
                 " poss2 " " " $(i, j-1)$
                 " poss3 " " " $(i-1, j)$

    return table(n,m) , (alignment)
              ↑
        overall cost

④ e.g.    x = STEP        y = APE

table

i↓
```
    A  P  E
S
T
E
P
```

```
    -  A  P  E
-   0  1  2  3
S   1  [1] ②
T   2
E   3
P   4
```

aux

```
         ins
    -  A  P  E
-      (0,0) (0,1) (0,2)
S  (0,0) (0,0) (0,1)
T  (1,0)
E  (2,0)
P  (3,0)
```

another aux table

```
    -  A  P  E
-   .  i  i  i
S   d  S  S
T   d
E   d
P   d
```

i=1 , j=1

subs  poss1 = table(0,0)+1 = 0+1 = 1
ins   poss2 = table(1,0)+1 = 1+1 = 2
del   poss3 = table(0,1)+1 = 1+1 = 2

i=1 , j=2

subs        table(0,1)+1  = 1+1 = 2
ins         table(1,1)+1  = 1+1 = 2
del         table(0,2)+1  = 2+1 = 3

```
- S
i
A
S
d  ?
-  A
```

```
- S
? s
A P
S -
? i
A P
```

Claim: Recurrence relation correctly computes min ⑤
cost.

Proof: By induction on $i$ and $j$.

Base case: $(i,j) = (0,0)$. In this case, we're aligning the empty string w/ itself, which has cost 0, and our base case of the recurrence correctly says $\text{cost}(0,0) = 0$. ✓

IH: $\text{cost}(i',j')$ is correct whenever $i' < i$ or $j' < j$.

Inductive step: Goal is to show $\text{cost}(i,j)$ is correct (assuming IH).

Consider the last op in an optimal alignment of $x[1..i]$ w/ $y[1..j]$.

Case 1: Last op is a subs.
   In this case, the cost of the alignment is
   $c(\text{subs}) + \underbrace{\text{cost of aligning } x[1..i-1] \,\&\, y[1..j-1]}$
                     $\text{by IH} = \text{cost}(i-1, j-1)$.
                     (applies b/c $i-1 < i$)
   $= \text{cost}(i-1, j-1) + c(\text{subs})$
   Because this was an opt alignment, this must be the min in the recurrence, so the recurrence relation correctly assigns
   $\text{cost}(i,j) = \text{cost}(i-1, j-1) + c(\text{subs})$.