

CS3000: Algorithms & Data

Jonathan Ullman

Lecture 7:

- Dynamic Programming: Segmented Least Squares

Sep 28, 2018

Dynamic Programming Recap

- **Recipe:**

difficult part

small

- (1) identify a set of **subproblems**
- (2) relate the subproblems via a **recurrence**

boilerplate

- (3) find an **efficient implementation** of the recurrence
- (4) **reconstruct the solution** from the DP table

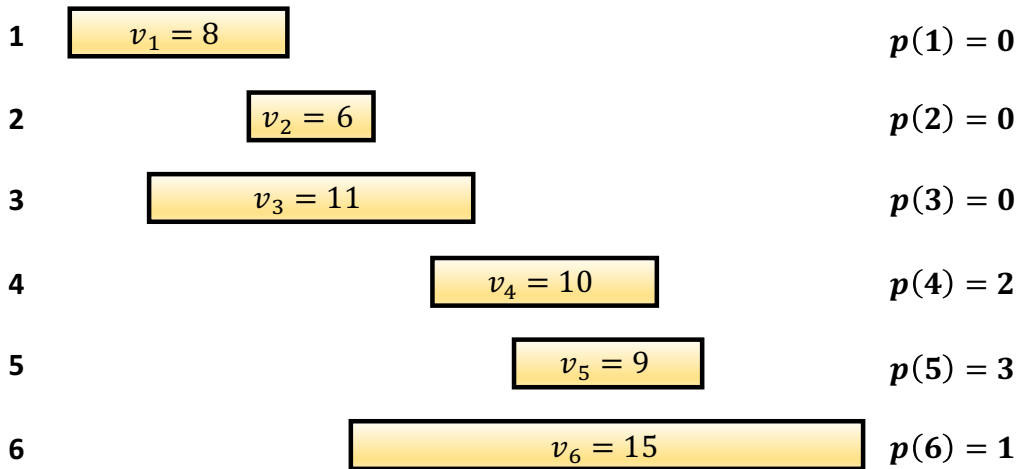
Dynamic Programming Recap

Goal: Find a schedule maximizing total value

$OPT(j) = \text{value of opt sched. for } \{1, \dots, j\}$

$$OPT(0) = 0$$

$$OPT(j) = \max \left\{ OPT(j-1), v_j + OPT(p(j)) \right\}$$



$$O_6 = \{6, 1\}$$

$$O_4 = \{4, 1\}$$

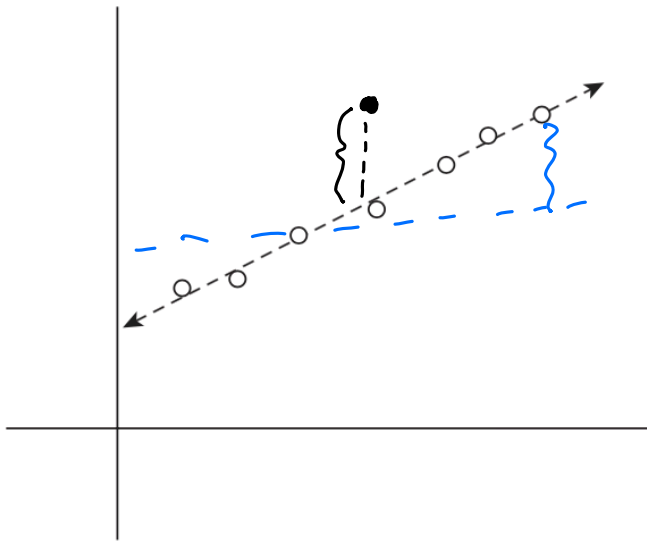
M[0]	M[1]	M[2]	M[3]	M[4]	M[5]	M[6]
0	8	8	11	18	20	23



Segmented Least Squares

Background: Least Squares

- **Input:** n data points $P = \{(x_1, y_1), \dots, (x_n, y_n)\}$
- **Output:** the line L (i.e. $y = ax + b$) that fits **best**
 - **best** = minimizes $error(L, P) = \sum_i (y_i - ax_i - b)^2$



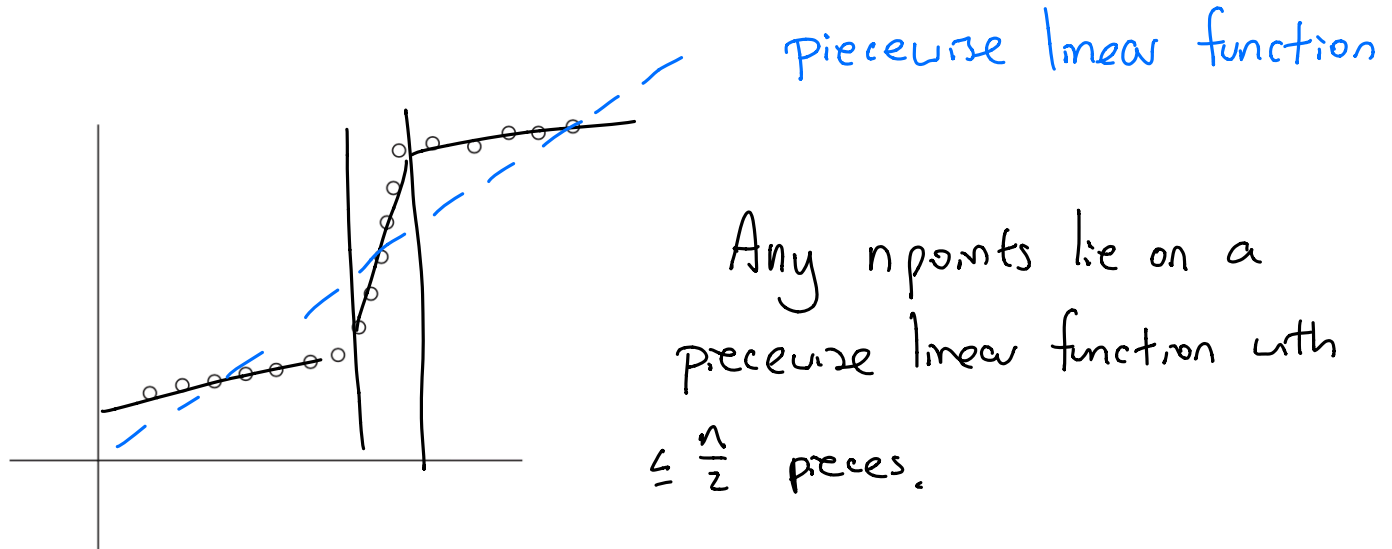
$$a = \frac{n\sum x_i y_i - (\sum x_i)(\sum y_i)}{n\sum x_i^2 - (\sum x_i)^2}$$

$$b = \frac{\sum y_i - a\sum x_i}{n}$$

Can find least squares estimate in $O(n)$ time

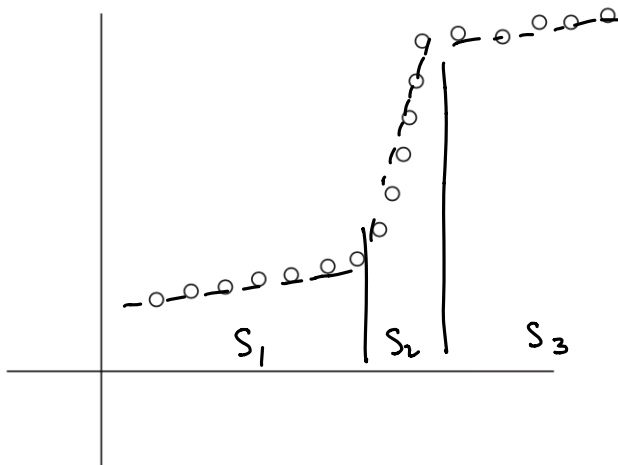
Segmented Least Squares

- **Input:** n data points $P = \{(x_1, y_1), \dots, (x_n, y_n)\}$
- What if the data does not look like a line?



Segmented Least Squares

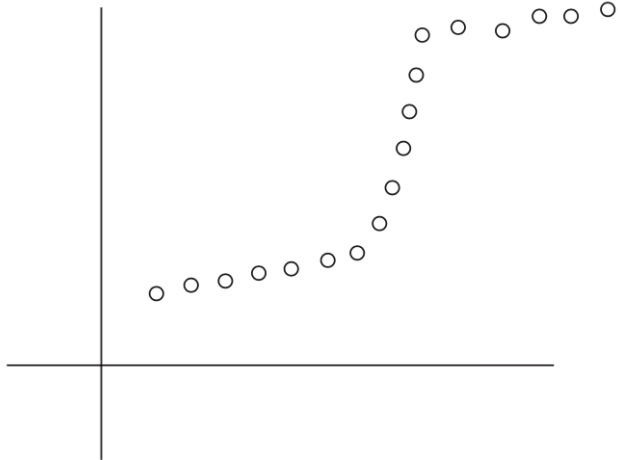
- **Input:** n data points $P = \{(x_1, y_1), \dots, (x_n, y_n)\}$, cost parameter $C > 0$
 - Assume $x_1 < x_2 < \dots < x_n$
- **Output:** a partition of P into contiguous segments S_1, S_2, \dots, S_m , lines L_1, L_2, \dots, L_m , minimizing “cost”



$$\begin{aligned} \text{cost}(S_1, \dots, S_m, L_1, \dots, L_m) \\ = C_m + \sum_{j=1}^m \text{error}(L_j, S_j) \end{aligned}$$

Segmented Least Squares

- **First observation:** for every segment S_j , L_j must be the (single) line of best fit for S_j
 - Let $L_{i,j}^*$ be the optimal line for $\{p_i, \dots, p_j\}$
 - Let $\varepsilon_{i,j} = \text{error}(L_{i,j}^*, \{p_i, \dots, p_j\})$



(can compute in $O(n^2)$ time.
(easy to do in $O(n^3)$)

SLS

Let $L_{i,j}^*$ be the optimal line for $\{p_i, \dots, p_j\}$

Let $\varepsilon_{i,j} = \text{error}(L_{i,j}^*, \{p_i, \dots, p_j\})$

- Let O be the **optimal** solution

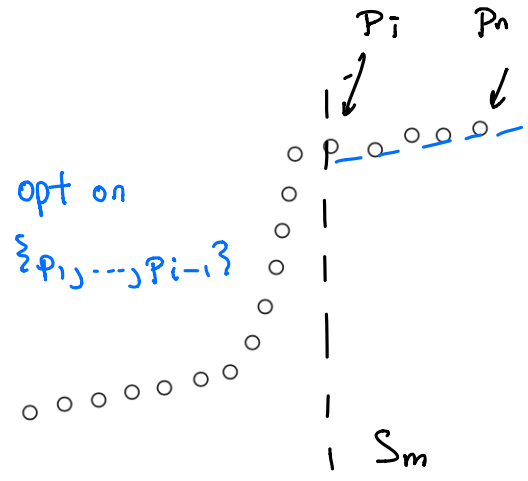
- O contains some final segment $\{p_i, \dots, p_n\}$ ($1 \leq i \leq n$)

- If the final segment is $\{p_i, \dots, p_n\}$
then the optimal segments are

① $\{p_i, \dots, p_n\}$

+ ② optimal segments for $\{p_1, \dots, p_{i-1}\}$

- cost would be $\varepsilon_{i,n} + C +$
opt cost on $\{p_1, \dots, p_{i-1}\}$



SLS

Let $L_{i,j}^*$ be the optimal line for $\{p_i, \dots, p_j\}$

Let $\varepsilon_{i,j} = \text{error}(L_{i,j}^*, \{p_i, \dots, p_j\})$

- Let $\text{OPT}(j)$ be the **value** of the optimal solution for points $\{p_1, \dots, p_j\}$ for $0 \leq j \leq n$
- **Case i :** final segment is $\{p_i, \dots, p_j\}$ for some $1 \leq i \leq j$
 - optimal solution is $L_{i,j}^* \cup$ optimal sol. for $\{p_1, \dots, p_{i-1}\}$
 - can use any $i \in \{1, \dots, j\}$

$$\text{OPT}(j) = \min_{i: 1 \leq i \leq j} \varepsilon_{i,j} + C + \text{OPT}(i-1)$$

$$\text{OPT}(0) = 0 \quad \text{OPT}(1) = C \quad \text{OPT}(2) = C$$

SLS

Let $L_{i,j}^*$ be the optimal line for $\{p_i, \dots, p_j\}$

Let $\varepsilon_{i,j} = \text{error}(L_{i,j}^*, \{p_i, \dots, p_j\})$

- Let $\text{OPT}(j)$ be the **value** of the optimal solution for points $\{p_1, \dots, p_j\}$
- **Case i :** final segment is $\{p_i, \dots, p_j\}$
 - optimal solution is $L_{i,j}^* \cup$ optimal sol. for $\{p_1, \dots, p_{i-1}\}$
 - can use any $i \in \{1, \dots, j\}$

Recurrence: $\text{OPT}(j) = \min_{1 \leq i \leq j} \varepsilon_{i,j} + C + \text{OPT}(i - 1)$

Base cases: $\text{OPT}(0) = 0$

$\text{OPT}(1) = \text{OPT}(2) = C$

SLS: Take I

```
// All inputs are global vars
FindOPT(n):
  if (n = 0): return 0
  elseif (n = 1,2): return C
  else:
    return  $\min_{1 \leq i \leq n} \epsilon_{i,n} + C + \text{FindOPT}(i - 1)$ 
```

$T(n) = \#$ of recursive calls

$$T(n) = \sum_{i=0}^{n-1} T(i) \gg T(n-1) + T(n-2)$$

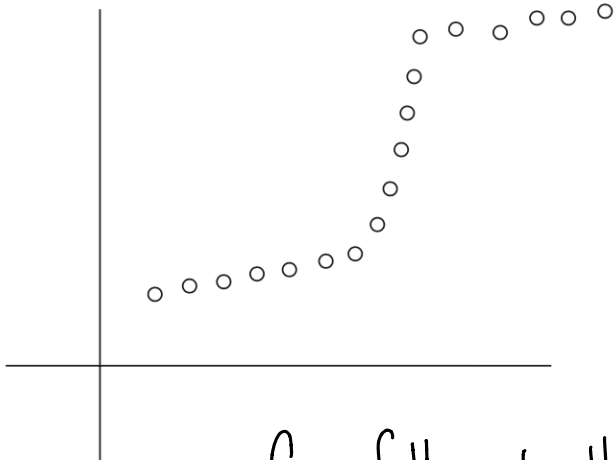
$$T(n) = \Omega(1.62^n)$$

SLS: Take II (“Top-Down”)

```
// All inputs are global vars
M ← empty array, M[0] ← 0, M[1] ← C, M[2] ← C
FindOPT(n):
  if (M[n] is not empty): return M[n]
  else:
    M[n] ←  $\min_{1 \leq i \leq n} \varepsilon_{i,n} + C + \text{FindOPT}(i - 1)$ 
    return M[n]
```

of recursive calls (n+1 arrays elements) × (≤ n calls)
= $O(n^2)$

SLS: Take III (“Bottom-Up”)

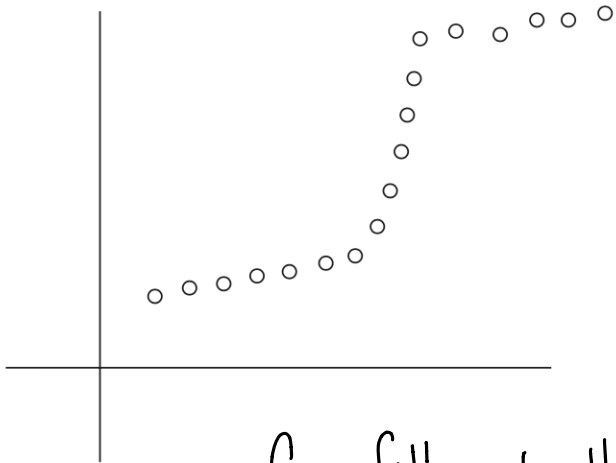


Can fill out the table left to right

M[0]	M[1]	M[2]	...	M[i]	...	M[n]
0	e	c	716.8



SLS: Take III (“Bottom-Up”)



Can fill out the table left to right

M[0]	M[1]	M[2]	...	M[i]	...	M[n]
0	e	c	716.8



SLS: Take III (“Bottom-Up”)

```
// All inputs are global vars
```

```
FindOPT(n):
```

```
  M[0] ← 0, M[1] ← C, M[2] ← C
```

```
  for (j = 3, ..., n):
```

$n-2$

```
    M[j] ←  $\min_{1 \leq i \leq j} \epsilon_{i,j} + C + M[i-1]$  ]  $O(n)$ 
```

```
  return M[n]
```

Running time is $O(n^2) + (\text{time to compute } \{\epsilon_{i,j}\})$

Rule of Thumb: running time $O(\# \text{ of dependencies btw problems})$

Finding Segments

Let $L_{i,j}^*$ be the optimal line for $\{p_i, \dots, p_j\}$

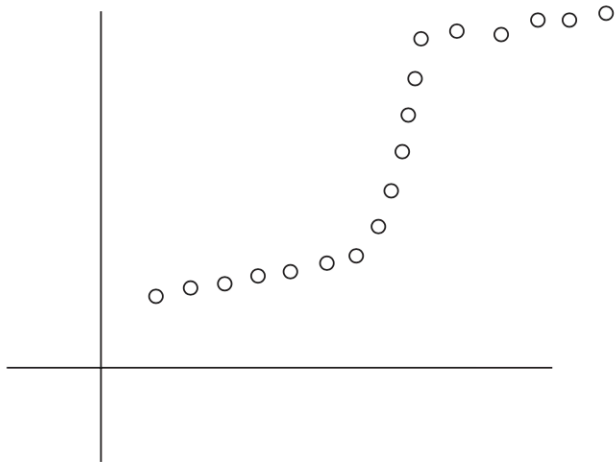
Let $\varepsilon_{i,j} = \text{error}(L_{i,j}^*, \{p_i, \dots, p_j\})$

- Let $\text{OPT}(j)$ be the **value** of the optimal solution for points $\{p_1, \dots, p_j\}$
- **Case i :** final segment is $\{p_i, \dots, p_j\}$
 - optimal solution is $L_{i,j}^* \cup$ optimal sol. for $\{p_1, \dots, p_{i-1}\}$
 - can use any $i \in \{1, \dots, j\}$
- If the final segment is $\{p_i, \dots, p_j\}$ then cost is $\varepsilon_{i,j} + C + \text{OPT}(i-1)$
- If $i \in \underset{1 \leq i \leq j}{\text{argmin}} \varepsilon_{i,j} + C + \text{OPT}(i-1)$ then there is an optimal solution uses $\{p_i, \dots, p_j\}$ as a final segment + optimal solution for $\{p_1, \dots, p_{i-1}\}$

Finding Segments

```
// All inputs are global vars
// M[0:n] contains solutions to subproblems
FindSol(M,n):
  if (n = 0): return  $\emptyset$ 
  elseif (n = 1): return {1}
  else:
    Let  $i \leftarrow \operatorname{argmax}_{1 \leq i \leq n} \epsilon_{i,n} + C + M[i - 1]$ :
    return {i, ..., n} + FindSol(M,i-1)
```

SLS: Take III (“Bottom-Up”)



$$M[n] = \epsilon_{i,j,n} + c + \text{OPT}(i-1)$$

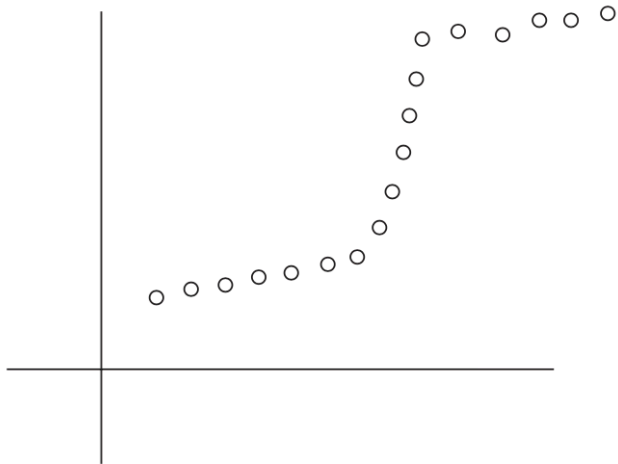
M[0]	M[1]	M[2]	...	M[i]	...	M[n]
0	c	c	716.8



Segmented Least Squares v.2

Segmented Least Squares v.2

- **Input:** n data points $P = \{(x_1, y_1), \dots, (x_n, y_n)\}$, parameter $1 \leq k \leq n$
 - Hard upper bound on the number of segments
- **Output:** a partition of P into $\leq k$ contiguous segments S_1, S_2, \dots, S_k minimizing “cost”



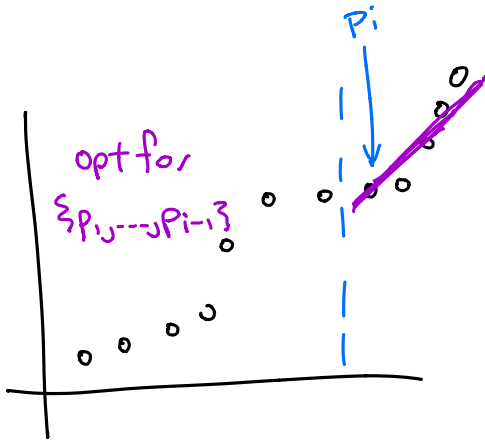
$$\begin{aligned} \text{cost}(S_1, \dots, S_k, L_1, \dots, L_k) \\ = \sum_{j=1}^k \text{error}(S_j, L_j) \end{aligned}$$

SLSv.2

Let $L_{i,j}^*$ be the optimal line for $\{p_i, \dots, p_j\}$

Let $\varepsilon_{i,j} = \text{error}(L_{i,j}^*, \{p_i, \dots, p_j\})$

- Let O be the optimal solution
- O uses some final segment $\{p_i, \dots, p_n\}$
- Possible recurrence:
$$\text{OPT}(n) = \min_{1 \leq i \leq n} \varepsilon_{i,n} + \text{OPT}(i-1)$$



problem: O_{i-1} uses k segments so we get $k+1$ segments total

idea: find the optimal solution using $k-1$ segments

SLSv.2

Let $L_{i,j}^*$ be the optimal line for $\{p_i, \dots, p_j\}$
Let $\varepsilon_{i,j} = \text{error}(L_{i,j}^*, \{p_i, \dots, p_j\})$

- Let $\text{OPT}(j, \ell)$ be the optimal solution for points $\{1, \dots, j\}$ using $\leq \ell$ segments
- **Case i :** final segment is $\{p_i, \dots, p_j\}$
 - optimal solution is $L_{i,j}^* \cup$ optimal solution for points $\{p_1, \dots, p_{i-1}\}$ using $\leq \ell - 1$ segments
 - can use any $i \in \{1, \dots, j\}$

Recurrence: $\text{OPT}(j, \ell) = \min_{1 \leq i \leq j} \varepsilon_{i,j} + \text{OPT}(i - 1, \ell - 1)$

Base cases: $\text{OPT}(0, \ell) = 0 \quad \forall \ell \geq 0$
 $\text{OPT}(j, 0) = \infty \quad \forall j \geq 1$

SLSv.2: Take II (“Top-Down”)

```
// All inputs are global vars
M ← empty array, M[0,ℓ] ← 0, M[j,0] ← ∞
FindOPT(n,k):
  if (M[n,k] is not empty): return M[n,k]
  else:
    M[n,k] ←  $\min_{1 \leq i \leq n} \varepsilon_{i,n} + \text{FindOPT}(i-1, k-1)$ 
    return M[n,k]
```

$(n+1)(k+1)$ subproblems

\times n calls per subproblem

$O(n^2k)$

SLSv.2: Take III ("Bottom-Up") $M[j, l] = \text{OPT}(j, l)$

Fill the table one column at a time
left to right

	$M[\cdot, 0]$	$M[\cdot, 1]$	$M[\cdot, 2]$	$M[\cdot, 3]$...	$M[\cdot, k]$
$M[0, \cdot]$	0	0	0	0	0	0
$M[1, \cdot]$	∞					
$M[2, \cdot]$	∞					
$M[3, \cdot]$	∞					
...	∞					
$M[n, \cdot]$	∞					

Because $\text{OPT}(j, l)$ depends on $\text{OPT}(i, l-1)$ for $i < j$ all "arrows" go up and left

SLSv.2: Take III (“Bottom-Up”)

```
// All inputs are global vars
```

```
FindOPT(n,k):
```

```
  M[0,ℓ] ← 0, M[j,0] ← ∞
```

```
  for (ℓ = 1, ..., k):
```

```
    for (j = 1, ..., n):
```

```
      M[j,ℓ] ←  $\min_{1 \leq i \leq j} \varepsilon_{i,j} + \text{FindOPT}(j-1, \ell-1)$  ]  $O(n)$ 
```

```
  return M[n,k]
```

$O(n^2k)$ time

$O(nk)$ space for M (+ $O(n^2)$ for $\{\varepsilon_{i,j}\}$)

SLSv.2: Finding Segments

```
// All inputs are global vars
// M[0:n,0:k] contains solutions to subproblems
FindSol(M,n,k):
  if (n = 0): return  $\emptyset$ 
  elseif (n = 1): return {1}
  else:
    let  $i \leftarrow \operatorname{argmax}_{1 \leq i \leq n} \varepsilon_{i,n} + M[i-1, k-1]$ :
    return {i, ..., n} + FindSol(M,i-1,k-1)
```

SLS Wrapup

- **Version 1:** can solve SLS with a “segment cost” in time $O(n^2)$ space $O(n^2)$
 - New idea: multiway case analysis for the final segment
- **Version 2:** can solve SLS with a “hard cap” of k segments in time $O(n^2k)$ space $O(n^2 + nk)$
 - New idea: introducing additional variables to expand the set of subproblems
- Correctness follows from the recurrence
- Computational costs:
 - Running time \approx total number of terms in all recurrences
 - Space \approx total number of subproblems