#### Segmented Least Squares



Given a collection P of points  $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n) \ldots$ 



Given a collection P of points  $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n) \ldots$ 

... find the "line  $\ell$  of best fit" through the points

 $\ell(x) = a \cdot x + b$ 



$$\ell(x) = a \cdot x + b \qquad \text{find } a \text{ and } b$$



What does "best" mean? Minimizes some error measure



**Residual:** difference between the *y*-coordinate of a point and the corresponding value of the fitted line



**Residual:** difference between the *y*-coordinate of a point and the corresponding value of the fitted line

**Error:** sum of the **squares** of the residuals

$$\mathsf{Err}(\ell, P) = \sum_{i=1}^{n} \left( \ell(x_i) - y_i \right)^2$$



The parameters of the line  $\ell(x) = a \cdot x + b$  of best fit are:

$$a = \frac{n \sum_{i} x_i y_i - \left(\sum_{i} x_i\right) \left(\sum_{i} y_i\right)}{n \sum_{i} x_i^2 - \left(\sum_{i} x_i\right)^2} \qquad b = \frac{\sum_{i} y_i - a \sum_{i} x_i}{n}$$



The parameters of the line  $\ell(x) = a \cdot x + b$  of best fit are:

$$a = \frac{n \sum_{i} x_i y_i - (\sum_{i} x_i) (\sum_{i} y_i)}{n \sum_{i} x_i^2 - (\sum_{i} x_i)^2} \qquad b = \frac{\sum_{i} y_i - a \sum_{i} x_i}{n}$$

Can be found in time O(n)

What if the points look like this?



What if the points look like this?



• No single line provides a good fit

What if the points look like this?



- No single line provides a good fit
- We can get a good fit with using **piecewise linear functions** instead of lines



- No single line provides a good fit
- We can get a good fit with using **piecewise linear functions** instead of lines

**Problem:** if we use piecewise linear functions, then we can trivially fit all points



- No single line provides a good fit
- We can get a good fit with using **piecewise linear functions** instead of lines

**Problem:** if we use piecewise linear functions, then we can trivially fit all points

**Idea:** each used segment incurs some  $\cot C$ 



- No single line provides a good fit
- We can get a good fit with using **piecewise linear functions** instead of lines

**Problem:** if we use piecewise linear functivially fit all points

Idea: each used segment incurs some  $\cot C$ 

Balances the quality of the fit with the # of used segments

## The Segmented Least Squares Problem

#### Input:

- A collection  $P = \{p_1 = (x_1, x_2), \dots, p_n = (x_n, y_n)\}$  of points with  $x_1 < x_2 < \dots < x_n$
- A cost  $C \in \mathbb{R}^+$

## The Segmented Least Squares Problem

#### Input:

- A collection  $P = \{p_1 = (x_1, x_2), \dots, p_n = (x_n, y_n)\}$  of points with  $x_1 < x_2 < \dots < x_n$
- A cost  $C \in \mathbb{R}^+$

#### Output:

- A partition of P into some number k of segments  $S_1, \ldots, S_k$
- where a segment is a subset of P containing a contiguous interval of points, i.e.,  $\{p_i, p_{i+1}, \ldots, p_j\}$  for some  $i \leq j$ ,
- that minimizes the *total penalty*  $C \cdot k + \sum_{i=1}^{k} \text{Err}(\ell_i, S_i)$ , where  $\ell_i$  is the best fit line for  $S_i$

## The Segmented Least Squares Problem

#### Input:

- A collection  $P = \{p_1 = (x_1, x_2), \dots, p_n = (x_n, y_n)\}$  of points with  $x_1 < x_2 < \dots < x_n$
- A cost  $C \in \mathbb{R}^+$

#### Output:

- A partition of P into some number k of segments  $S_1, \ldots, S_k$
- where a segment is a subset of P containing a contiguous interval of points, i.e.,  $\{p_i, p_{i+1}, \ldots, p_j\}$  for some  $i \leq j$ ,
- that minimizes the *total penalty*  $C \cdot k + \sum_{i=1}^{k} \text{Err}(\ell_i, S_i)$ , where  $\ell_i$  is the best fit line for  $S_i$   $\text{Err}(S_i)$

## Example (qualitative)



#### C = 5

# Example (qualitative)



**One segment:** Total penalty:  $1 \cdot 5 + 200 = 205$ 



**One segment:** Total penalty:  $1 \cdot 5 + 200 = 205$ **Two segments:** Total penalty:  $2 \cdot 5 + 50 = 60$ 

# Example (qualitative) $\ell_3$ C=5

**One segment:** Total penalty:  $1 \cdot 5 + 200 = 205$ **Two segments:** Total penalty:  $2 \cdot 5 + 50 = 60$ **Three segments:** Total penalty:  $3 \cdot 5 + 40 = 55$ 



One segment: Total penalty:  $1 \cdot 5 + 200 = 205$ Two segments: Total penalty:  $2 \cdot 5 + 50 = 60$ Three segments: Total penalty:  $3 \cdot 5 + 40 = 55$ Four segments: Total penalty:  $4 \cdot 5 + 37 = 57$ 

# Example (qualitative)



One segment: Total penalty:  $1 \cdot 5 + 200 = 205$ Two segments: Total penalty:  $2 \cdot 5 + 50 = 60$ Three segments: Total penalty:  $3 \cdot 5 + 40 = 55$ Four segments: Total penalty:  $4 \cdot 5 + 37 = 57$ 

Let OPT(i) denote the penalty incurred by an optimal solution for the instance that only considers the points in  $\{p_1, p_2, \ldots, p_i\}$ 

Let OPT(i) denote the penalty incurred by an optimal solution for the instance that only considers the points in  $\{p_1, p_2, \ldots, p_i\}$ 

If we know that the last segment in an optimal solution is  $p_j, \ldots, p_n$ , then the value of the optimal solution is

 $\mathsf{OPT}(n) = C + \mathsf{Err}(\{p_j, \dots, p_n\}) + \mathsf{OPT}(j-1).$ 

Let OPT(i) denote the penalty incurred by an optimal solution for the instance that only considers the points in  $\{p_1, p_2, \ldots, p_i\}$ 

If we know that the last segment in an optimal solution is  $p_j, \ldots, p_n$ , then the value of the optimal solution is

 $\mathsf{OPT}(n) = C + \mathsf{Err}(\{p_j, \dots, p_n\}) + \mathsf{OPT}(j-1).$ 



Let OPT(i) denote the penalty incurred by an optimal solution for the instance that only considers the points in  $\{p_1, p_2, \ldots, p_i\}$ 

If we know that the last segment in an optimal solution is  $p_j, \ldots, p_n$ , then the value of the optimal solution is

 $\mathsf{OPT}(n) = C + \mathsf{Err}(\{p_j, \dots, p_n\}) + \mathsf{OPT}(j-1).$ 



Let OPT(i) denote the penalty incurred by an optimal solution for the instance that only considers the points in  $\{p_1, p_2, \ldots, p_i\}$ 

If we know that the last segment in an optimal solution is  $p_j, \ldots, p_n$ , then the value of the optimal solution is

 $\mathsf{OPT}(n) = C + \mathsf{Err}(\{p_j, \dots, p_n\}) + \mathsf{OPT}(j-1).$ 



Let OPT(i) denote the penalty incurred by an optimal solution for the instance that only considers the points in  $\{p_1, p_2, \ldots, p_i\}$ 

If we know that the last segment in an optimal solution is  $p_j, \ldots, p_n$ , then the value of the optimal solution is

 $\mathsf{OPT}(n) = C + \mathsf{Err}(\{p_j, \dots, p_n\}) + \mathsf{OPT}(j-1).$ 



Let OPT(i) denote the penalty incurred by an optimal solution for the instance that only considers the points in  $\{p_1, p_2, \ldots, p_i\}$ 

If we know that the last segment in an optimal solution is  $p_j, \ldots, p_n$ , then the value of the optimal solution is

 $\mathsf{OPT}(n) = C + \mathsf{Err}(\{p_j, \dots, p_n\}) + \mathsf{OPT}(j-1).$ 



Let OPT(i) denote the penalty incurred by an optimal solution for the instance that only considers the points in  $\{p_1, p_2, \ldots, p_i\}$ 

If we know that the last segment in an optimal solution is  $p_j, \ldots, p_n$ , then the value of the optimal solution is

 $\mathsf{OPT}(n) = C + \mathsf{Err}(\{p_j, \dots, p_n\}) + \mathsf{OPT}(j-1).$ 



Let OPT(i) denote the penalty incurred by an optimal solution for the instance that only considers the points in  $\{p_1, p_2, \ldots, p_i\}$ 

If we know that the last segment in an optimal solution is  $p_j, \ldots, p_n$ , then the value of the optimal solution is

 $\mathsf{OPT}(n) = C + \mathsf{Err}(\{p_j, \dots, p_n\}) + \mathsf{OPT}(j-1).$ 



Let OPT(i) denote the penalty incurred by an optimal solution for the instance that only considers the points in  $\{p_1, p_2, \ldots, p_i\}$ 

If we know that the last segment in an optimal solution is  $p_j, \ldots, p_n$ , then the value of the optimal solution is

 $\mathsf{OPT}(n) = C + \mathsf{Err}(\{p_j, \dots, p_n\}) + \mathsf{OPT}(j-1).$ 



# A Dynamic Programming Algorithm

#### Subproblem definition:

 $OPT(i) = penalty incurred by an optimal solution for the instance that only considers the points in <math>\{p_1, p_2, \dots, p_i\}$ Base case: OPT(0) = 0

Recursive formula:

```
(for i \geq 1)
```

$$\mathsf{OPT}(i) = \min_{j=1,...,i} \left\{ C + \mathsf{Err}(\{p_j,\ldots,p_n\}) + \mathsf{OPT}(j-1) \right\}.$$

**Order of subproblems:** OPT(1), OPT(2), ..., OPT(n)

**Solution:** OPT(n)

How many subproblems?

How many subproblems?



How many subproblems?



#### How much time per subproblem?

- We need to test O(n) choices of j.
- How much time per choice of *j*?

How many subproblems?



#### How much time per subproblem?

- We need to test O(n) choices of j.
- How much time per choice of j?
  - Need to find the best fit line  $\ell$  for  $\{p_j, \ldots, p_i\}$  O(n)
  - Need to find the error  $\operatorname{Err}(\ell, \{p_j, \ldots, p_i\})$  O(n)

O(n)

How many subproblems?

How much time per subproblem?  $O(n^2)$ 

- We need to test O(n) choices of j.
- How much time per choice of j?
  - Need to find the best fit line  $\ell$  for  $\{p_j, \ldots, p_i\}$  O(n)
  - Need to find the error  $\operatorname{Err}(\ell, \{p_j, \ldots, p_i\})$  O(n)

**Overall time:**  $O(n^3)$ 

O(n)

How many subproblems?

How much time per subproblem?  $O(n^2)$ 

- We need to test O(n) choices of j.
- How much time per choice of j?
  - Need to find the best fit line  $\ell$  for  $\{p_j, \ldots, p_i\}$  O(n)
  - Need to find the error  $\operatorname{Err}(\ell, \{p_j, \ldots, p_i\})$  O(n)

**Overall time:**  $O(n^3)$ 

Can we do better?

#### Improving the Time Complexity

**Goal:** find the best fit line  $\ell(x) = a \cdot x + b$ , and the error **Err** $(\ell, \{p_j, \dots, p_i\})$  **quickly** 

$$a = \frac{n \sum_{h=j}^{i} x_h y_h - \left(\sum_{h=j}^{i} x_h\right) \left(\sum_{h=j}^{i} y_h\right)}{n \sum_{h=j}^{i} x_h^2 - \left(\sum_{h=j}^{i} x_h\right)^2}$$

$$b = \frac{\sum_{h=j}^{i} y_h - a \sum_{h=j}^{i} x_h}{n}$$

$$\mathsf{Err}(\ell, \{p_j, \dots, p_i\}) = \sum_{h=j}^i \left( \ell(x_h) - y_h \right)^2$$





Can we find this quantity quickly?





Can we find this quantity quickly?



Can we find this quantity quickly?



$$b = \frac{\sum_{h=j}^{i} y_h}{n} - a \sum_{h=j}^{i} x_h}{n}$$

Can we find this quantity quickly?

What about this one?

s one?  

$$n\sum_{h=j}^{i} x_{h}y_{h} - \left(\sum_{h=j}^{i} x_{h}\right) \left(\sum_{h=j}^{i} y_{h}\right)$$
What about  

$$a = \frac{1}{n\sum_{h=j}^{i} x_{h}^{2}} - \left(\sum_{h=j}^{i} x_{h}\right)^{2}$$
What about  
this one?

$$b = \frac{\sum_{h=j}^{i} y_h - a \sum_{h=j}^{i} x_h}{n}$$

Can we find this quantity quickly?

What about this one?

s one?  

$$n \sum_{h=j}^{i} x_{h} y_{h} - \left(\sum_{h=j}^{i} x_{h}\right) \left(\sum_{h=j}^{i} y_{h}\right) \text{What about this one?}$$

$$a = \frac{n \sum_{h=j}^{i} x_{h}^{2} - \left(\sum_{h=j}^{i} x_{h}\right)^{2}}{n \sum_{h=j}^{i} x_{h}^{2} - \left(\sum_{h=j}^{i} x_{h}\right)^{2}}$$
What about this one?  

$$b = \frac{\sum_{h=j}^{i} y_{h} - a \sum_{h=j}^{i} x_{h}}{n}$$

All the marked quantities can be found in **constant time** after a linear-time preprocessing using **prefix sums**.

$$\mathsf{Err}(\ell, \{p_j, \dots, p_i\}) = \sum_{h=j}^{i} \left(\ell(x_h) - y_h\right)^2 \qquad \ell(x) = a \cdot x + b$$

$$\mathsf{Err}(\ell, \{p_j, \dots, p_i\}) = \sum_{h=j}^{i} \left(\ell(x_h) - y_h\right)^2 \qquad \ell(x) = a \cdot x + b$$

$$= \sum_{h=j}^{i} \left( \ell(x_h)^2 + y_h^2 - 2\ell(x_h)y_h \right)$$

$$\mathsf{Err}(\ell, \{p_j, \dots, p_i\}) = \sum_{h=j}^{i} \left(\ell(x_h) - y_h\right)^2 \qquad \qquad \ell(x) = a \cdot x + b$$

$$=\sum_{h=j}^{i} \left( a^2 x_h^2 + b^2 + 2ax_h b + y_h^2 - 2ax_h y_h - 2by_h \right)$$

$$\mathsf{Err}(\ell, \{p_j, \dots, p_i\}) = \sum_{h=j}^{i} \left(\ell(x_h) - y_h\right)^2 \qquad \qquad \ell(x) = a \cdot x + b$$

$$=a^{2}\sum_{h=j}^{i}x_{h}^{2}+(j-i+1)b^{2}+2ab\sum_{h=j}^{i}x_{h}+\sum_{h=j}^{i}y_{h}^{2}$$
$$-2a\sum_{h=j}^{i}x_{h}y_{h}-2b\sum_{h=j}^{i}y_{h}$$

$$\mathsf{Err}(\ell, \{p_j, \dots, p_i\}) = \sum_{h=j}^{i} \left(\ell(x_h) - y_h\right)^2 \qquad \qquad \ell(x) = a \cdot x + b$$

$$=a^{2}\sum_{h=j}^{i}x_{h}^{2}+(j-i+1)b^{2}+2ab\sum_{h=j}^{i}x_{h}+\sum_{h=j}^{i}y_{h}^{2}$$
$$-2a\sum_{h=j}^{i}x_{h}y_{h}-2b\sum_{h=j}^{i}y_{h}$$

$$\mathsf{Err}(\ell, \{p_j, \dots, p_i\}) = \sum_{h=j}^{i} \left(\ell(x_h) - y_h\right)^2 \qquad \qquad \ell(x) = a \cdot x + b$$

$$=a^{2}\sum_{h=j}^{i}x_{h}^{2}+(j-i+1)b^{2}+2ab\sum_{h=j}^{i}x_{h}+\sum_{h=j}^{i}y_{h}^{2}$$
$$-2a\sum_{h=j}^{i}x_{h}y_{h}-2b\sum_{h=j}^{i}y_{h}$$

$$\mathsf{Err}(\ell, \{p_j, \dots, p_i\}) = \sum_{h=j}^{i} \left(\ell(x_h) - y_h\right)^2 \qquad \qquad \ell(x) = a \cdot x + b$$

$$=a^{2}\sum_{h=j}^{i}x_{h}^{2} + (j-i+1)b^{2} + 2ab\sum_{h=j}^{i}x_{h} + \sum_{h=j}^{i}y_{h}^{2}$$
$$-2a\sum_{h=j}^{i}x_{h}y_{h} - 2b\sum_{h=j}^{i}y_{h}$$

# Finding $\operatorname{Err}(\cdot, \cdot)$ quickly

$$\mathsf{Err}(\ell, \{p_j, \dots, p_i\}) = \sum_{h=j}^{i} \left(\ell(x_h) - y_h\right)^2 \qquad \qquad \ell(x) = a \cdot x + b$$

$$=a^{2}\sum_{h=j}^{i}x_{h}^{2} + (j-i+1)b^{2} + 2ab\sum_{h=j}^{i}x_{h} + \sum_{h=j}^{i}y_{h}^{2}$$
$$-2a\sum_{h=j}^{i}x_{h}y_{h} - 2b\sum_{h=j}^{i}y_{h}$$

$$\mathsf{Err}(\ell, \{p_j, \dots, p_i\}) = \sum_{h=j}^{i} \left(\ell(x_h) - y_h\right)^2 \qquad \qquad \ell(x) = a \cdot x + b$$

$$= a^{2} \sum_{h=j}^{i} x_{h}^{2} + (j - i + 1)b^{2} + 2ab \sum_{h=j}^{i} x_{h} + \sum_{h=j}^{i} y_{h}^{2}$$
$$- 2a \sum_{h=j}^{i} x_{h}y_{h} - 2b \sum_{h=j}^{i} y_{h}$$

All the marked quantities can be found in **constant time** after a linear-time preprocessing using **prefix sums**.

O(n)

How many subproblems?

How much time per subproblem?  $O(n^2)$ 

- We need to test O(n) choices of j.
- How much time per choice of *j*?
  - Need to find the best fit line  $\ell$  for  $\{p_j, \ldots, p_i\}$  O(n)

- Need to find the error  $Err(\ell, \{p_j, \ldots, p_i\})$  O(n)

How many subproblems?

How much time per subproblem?

 $O(n) \mathcal{O}$ 

O(n)

- We need to test O(n) choices of j.
- How much time per choice of *j*?
  - Need to find the best fit line  $\ell$  for  $\{p_j, \ldots, p_i\}$
  - Need to find the error  $Err(\ell, \{p_j, \ldots, p_i\})$

How many subproblems?

How much time per subproblem?

O(n)

O(n)

- We need to test O(n) choices of j.
- How much time per choice of *j*?
  - Need to find the best fit line  $\ell$  for  $\{p_j, \ldots, p_i\}$
  - Need to find the error  $Err(\ell, \{p_j, \ldots, p_i\})$

+O(n)-time preprocessing (one-time only)

O(n)

O(n)

How many subproblems?

How much time per subproblem?

- We need to test O(n) choices of j.
- How much time per choice of *j*?
  - Need to find the best fit line  $\ell$  for  $\{p_j, \ldots, p_i\}$
  - Need to find the error  $Err(\ell, \{p_j, \ldots, p_i\})$

+O(n)-time preprocessing (one-time only)

**Overall time:**  $O(n^2)$