

Main idea :-
Line Fitting basics :-

Set P :- n points $\Rightarrow (x_1, y_1) - \dots - (x_n, y_n)$

Suppose $x_1 < x_2 < \dots < x_n$

Line L :- $y = ax + b$

Error of line L with respect to P :-

\Rightarrow sum of squared distances to the point P

$$\text{Error}(L, P) = \sum_{i=1}^n (y_i - ax_i - b)^2$$



Gives better approximation of how well line L fits P points.

less error \Rightarrow better fit

$\exists P$ $n \geq 2$,

Given set P points

The line with minimal error $\text{Error}(L, P)$

$$\text{is } y = ax + b$$

$$a = \frac{n \sum_{i=1}^n x_i y_i - \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right)}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2}$$

$$b = \frac{\sum_{i=1}^n y_i - a \sum_{i=1}^n x_i}{n}$$

Multiline fitting problem :-

Here, we fit points using multiple lines.

→ not simple to know the number of lines that are best choice beforehand.

Solve it as a modelling problem :-

- model a set of observed points $\{x_i, y_i\}$ using a simple & accurate model

Here :-

model \Rightarrow set of lines, each with different x interval

Input :-

① Set of n -points in a 2-D plane.

$$P = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

$$x_1 < x_2 < \dots < x_n$$

$$P_i = (x_i, y_i) \quad \text{for } i = 1 \text{ to } n$$

② Real number $C > 0$

Output :-

Partition of points P into k segments :-

$$P_1 = \{P_1, P_2, \dots, P_{m_1}\}$$

$$P_2 = \{P_{m_1+1}, P_{m_1+2}, \dots, P_{m_2}\}$$

$$P_3 = \{P_{m_2+1}, P_{m_2+2}, \dots, P_{m_3}\}$$

⋮

$$P_n = \{ P_{m_{k-1}+1}, \dots, P_{m_k} = P_n \}$$

then minimize the cost function

$$\text{cost}(P, C) = \underbrace{\epsilon(P_1) + \epsilon(P_2) + \dots + \epsilon(P_k)}_{\text{total error using } k \text{ lines to fit } n \text{- points}} + C_k,$$

\downarrow
total error using
k lines to fit
n-points

\downarrow
penalty
for using k
lines

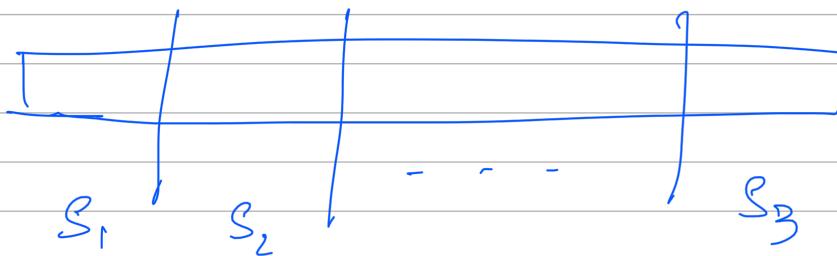
The value of k & turning points $(P_{m_1}, P_{m_2}, \dots, P_{m_{k-1}})$
are needed to be decided by the algorithm

Idea behind my algorithm:-

I have solved this problem using dynamic programming.

I took the inspiration from the rod cutting problem.

In the given problem, we need to segregate the points in such a way that, the number of lines chosen and the cost has to be optimal.



In the rod cutting problem we go for a solution that gives maximum price. But, in this problem we should go for a solution that gives minimal cost.

The optimal. cost is given by

$d[i] \Rightarrow$ cost till i^{th} point

{ if $d_p[i] > 0$

$$d_p[i-k] + \text{cost}[i-k+1][i]$$

$$d_p[i] = \min$$

else

$$\text{cost}[0][i]$$

$$2 \leq k \leq i$$



because we cannot abandon the last point, we must always include it

Here k is the number of steps we are going back from the current dp point to get the cost.

In order to get the k list, we keep track of break points every time we find a better cost dp from the point that we are standing and add it to the breakpoints of the current dp's corresponding index.

K Values are stored such that

$$K[i] = [i] + \begin{cases} \text{if } :- \\ (dp[i-K] + cost[i-K+1][i]) \leftarrow \min \\ K \Rightarrow K[i-K] \\ \text{else} \\ K \Rightarrow [] \end{cases}$$

we just have to keep track of the optimal dp's break points and append them to the current dp along with the end point till the current dp

Here $dp[i] \Rightarrow$ minimum cost till i^{th} point
 $0 \leq i \leq \text{numpoints}$

$K[i] \Rightarrow$ Break points till i^{th} dp point
 $0 \leq i \leq \text{numpoints}$

In this algorithm I calculate cost between two points using an error matrix, which uses $O(n^2)$ to get the entire error matrix

$cost[i][j] \Rightarrow$ minimum cost from point i to j .

Pseudo Code :-

Main function to pass the dataset :-

piecewise dict(dataset)

For i from (0 to len of dataset) :-

// read the data set to first record
cur_dataset = dataset[i]

cur_dataset сент = piecewise_linear_segmentation
(cur_dataset)

Final result.append (cur_dataset сент)

return Final_result

Calculate сент for each record :-



Piecewise linear Segmentation (record)

load n_list, x_list_y_list, c_list from record

// error matrix

errors = linear_regression_error_matrix(n_list, x_list, y_list, c_list)

// errors[i][j] = cost for segment from ith to jth point

dp = [0, 0, ..., 0] // size n_list - 1

dp[i] = errors[0][i]

K = [[], [], ..., []] // size n_list - 1

K[i].append(i)

for i from 2 to (n_list - 1)

min_val = MAX_INT

K[i].append(i)

tempList = []

for k from 2 to i // steps to go back

if (i-k) > 0 :

$\text{temp_val} = \text{dp}[i-k] + \text{error_list}[i-k+1][i]$
 if ($\text{temp_val} < \text{min_val}$):
 $\text{min_val} = \text{temp_val}$
 $\text{temp_list} = K[i-k]$

Else :

$\text{temp_val} = \text{error_list}[0][i]$
 if ($\text{temp_val} < \text{min_val}$):
 $\text{min_val} = \text{temp_val}$
 $\text{temp_list} = []$

$\text{dp}[i] = \text{min_val}$ // update dp
 $K[i] = K[i] + \text{temp_list}$ // update K

$\text{result}['K_list'] = \text{len}(K[n_list-1])$
 $\text{result}['\text{last_points_list}'] = K[n_list-1]$
 $\text{result}['\text{OPT_list}'] = \text{dp}[n_list-1]$

Return result.

✓ calculate error matrix :-

Linear_regression_error_matrix(n_list, x_list, y_list, c_list)

error_matrix = np.zeros((n_list, n_list))

// get the cumulative sum arrays // cumulative sum till point

$x_sum, y_sum, x^2_sum, y^2_sum, xy_dot$
 $= \text{preprocess_arrays}(n_list, x_list, y_list)$

for i from 0 to n_list :

for j from i to n_list :

$n = j - i + 1$

if ($j - i = 1$):

$\text{error_matrix}[i][j] = c_{list}$

if ($i == 0$):

$$X_{sum_ij} = X_{sum}[j]$$
$$Y_{sum_ij} = Y_{sum}[j]$$

$$X_{pow_2_sum_ij} = X_{pow_2_sum}[j]$$

$$Y_{pow_2_sum_ij} = Y_{pow_2_sum}[j]$$

$$XY_dot_ij = XY_dot[j]$$

else:

$$X_{sum_ij} = X_{sum}[j] - X_{sum}[i-1]$$

$$Y_{sum_ij} = Y_{sum}[j] - Y_{sum}[i-1]$$

$$X_{pow_2_sum_ij} = X_{pow_2_sum}[j] - X_{pow_2_sum}[i-1]$$

$$Y_{pow_2_sum_ij} = Y_{pow_2_sum}[j] - Y_{pow_2_sum}[i-1]$$

$$XY_dot_ij = XY_dot[j] - XY_dot[i-1]$$

$$a = ((n * XY_dot_ij) - (X_{sum_ij} * Y_{sum_ij}))$$

$$((n * X_{pow_2_sum_ij}) - (X_{sum_ij} * Y_{sum_ij}))$$

$$b = (Y_{sum_ij} - a * X_{sum_ij}) / (n)$$

error_matrix[i][j] = absolute value of

$$\sqrt{(Y_{pow_2_sum_ij})^2 + (a^2 * X_{pow_2_sum_ij})^2 + (n * b^2) + (2 * ab * X_{sum_ij}) - (2 * a * XY_dot_ij) - (2 * b * Y_{sum_ij}) + c \cdot \text{list}}$$

✓
code to find cummulative sum of variables -

preprocess arrays (x_list , y_list , x^y_list):

$$x_sum[0] = x_list[0], y_sum[0] = y_list[0],$$
$$x_pow_2_sum[0] = x_list[0] * x_list[0],$$
$$y_pow_2_sum[0] = y_list[0] * y_list[0]$$
$$xy_dot[0] = x_list[0] * y_list[0]$$

// calculate cumulative sum for other points

for i in range($0, n_list$): -

$$x_sum[i] = x_sum[i-1] + x_list[i]$$
$$y_sum[i] = y_sum[i-1] + y_list[i]$$
$$x_pow_2_sum[i] = x_pow_2_sum[i-1] + (x_list[i] * x_list[i])$$
$$y_pow_2_sum[i] = y_pow_2_sum[i-1] + (y_list[i] * y_list[i])$$
$$xy_dot[i] = xy_dot[i-1] + (x_list[i] * y_list[i])$$

return $x_sum, y_sum, x_pow_2_sum, y_pow_2_sum, xy_dot$

Proof of correctness

My algorithm is correct because it covers various cases in which no single point is abandoned & always considers two or more points in a segment.

In my algorithm, I collect and save the cumulative cost till a point i in an array called dp .

$dp[0] \Rightarrow$ cost till 0^{th} point but this is considered because single point can't be abandoned.

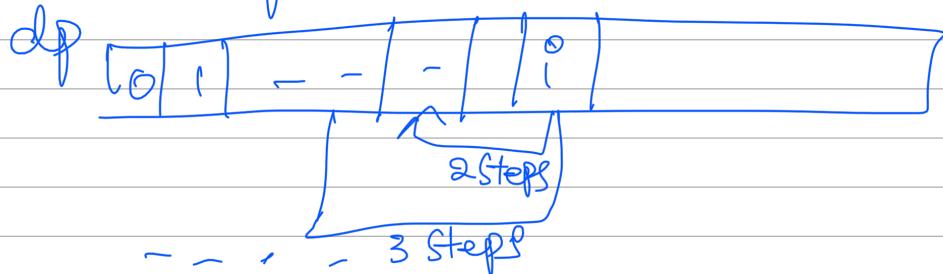
$dp[1] \Rightarrow$ cost $[0 \rightarrow 1]$

whenever we encounter $dp[0]$ or $dp[i]$

↓
we always calculate the cost
as $\text{error}(0 \rightarrow i)$

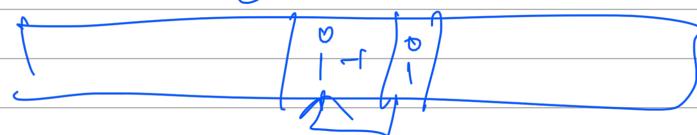
as we cannot abandon the 0th point.

When we are at point 'i'



we always take 2 or more steps back
to find the optimal cost.

If we take just one step:-



it will be $dp[i-1] + \text{cost}(i)$

↓
single point

So we do not consider the cost
with single step.

For every iteration we check for all
possible cases & find the least

$dp[0] \Rightarrow \text{NA}$

$dp[i] \Rightarrow e[0 \rightarrow i]$

$$dp[2] \Rightarrow \cancel{dp[0]} + e[0 \rightarrow 2]$$

2 steps back

$$+ e[1 \rightarrow 2]$$

$$+ e[0 \rightarrow 2]$$

as we encountered $0 \rightarrow i$ in dp index we only consider
 $e[0 \rightarrow i]$ whenever we encounter $dp[0]$ or $dp[-1]$

$$dp[3] \Rightarrow \min \left(\begin{array}{l} dp[i] + e[2 \rightarrow 3] \\ dp[0] + e[1 \rightarrow 3] \\ dp[-1] + e[0 \rightarrow 3] \end{array} \right)$$

2 step 3 step 4 steps

$$dp[u] = \min \left(\begin{array}{l} dp[2] + e[3 \rightarrow u] \\ dp[i] + e[2 \rightarrow u] \\ dp[0] + e[1 \rightarrow u] \\ dp[-1] + e[0 \rightarrow u] \end{array} \right)$$

2 steps 3 steps 4 steps 5 steps

Every time we update the minimum, we update 'k' value corresponding to the considered dp that gives the minimum

As my algorithm covers all possible cases when solving a problem such as multiline fitting I can say that this is an optimal approach with the best time complexity $O(n^2)$ and this approach also involves dynamic programming.

Time Complexity analysis

The time complexity for my algorithm is $O(n^2)$ for a dataset record (single set of points)

The following is the flow for my algorithm :-

```
In [5]: def piecewise_linear_segmentation(record):
    n_list = record[0]
    x_list = record[1]
    y_list = record[2]
    c_list = record[3]

    # calculate and store errors for the possible ranges
    errors = linear_regressor_error_matrix(n_list,x_list,y_list,c_list)

    # Lets Calculate (k_list, last_points_list, OPT_list)
    # k_list = number of segments
    # last_points_list = list of last points in each segment
    # OPT_list is the cost function value

    # Lets create a list dp to store the cost till a point (this stores the OPT_list value)
    # dp[i] means cost till ith point
    dp = [0 for i in range(n_list)] → O(1)
    dp[1] = errors[0][1] → O(1)

    # track K value for every i value in DP (this stores the last_points_list in K)
    K = [[] for i in range(n_list)] → O(n)
    K[1].append(1) → O(1)

    # run for ith point cost in dp
    for i in range(2,n_list):
        #store min for this dp point
        min_val = sys.maxsize → O(1)
        K[i].append(i)
        temp_list = []
        #number of steps we are going back to check the cost
        for k in range(2,i+1):
            if ((i-k) > 0):
                temp_val = dp[i-k] + errors[i-k+1][i] → O(1)
                if(temp_val < min_val):
                    min_val = temp_val
                    temp_list = K[i-k] → O(1)
            else:
                temp_val = errors[0][i]
                if(temp_val < min_val):
                    min_val = temp_val
                    temp_list = []
            min_val = min(temp_val,min_val) → O(1)
            temp_list = []

        dp[i] = min_val
        K[i] = K[i]+temp_list → O(1)

    print("K list is")
    print(K[n_list-1])
    print("Optimum value is")
    print(dp[n_list-1])
    print("Segments count is")
    print(len(K[n_list-1]))
    result = {}
    result['k_list'] = len(K[n_list-1])
    result['last_points_list'] = K[n_list-1]
    result['OPT_list'] = dp[n_list-1] → O(1)
    return result
```

Total ($O(n^2)$)

$\rightarrow O(n^2)$

$\rightarrow O(n^2)$

Total = $O(n^2)$

Calculate Error Matrix

```
In [4]: def linear_regressor_error_matrix(n_list,x_list,y_list,c_list):
    #create an error matrix for the given points
    #error_matrix[i][j] => error from point i to j segment
    #i from 0 to n_list-1 and j from 0 to n_list-1
    error_matrix = np.zeros((n_list, n_list)) - O(1)
    # let's do the cumulative sum parameters
    X_sum,Y_sum,X_pow_2_sum,Y_pow_2_sum,XY_dot = preprocess_arrays(n_list,x_list,y_list)
    for i in range(0,n_list):
        for j in range(i+1,n_list):
            if(i==j):
                error_matrix[i][j] = c_list[j]
            elif(i==0):
                #let's write the if condition where it checks if 0th point is involved
                #here we can directly use cumulative sum till jth point
                total_points = n
                n = j+1
                #calc individual params
                X_sum_ij = X_sum[0]
                Y_sum_ij = Y_sum[0]
                X_pow_2_sum_ij = X_pow_2_sum[0]
                Y_pow_2_sum_ij = Y_pow_2_sum[0]
                XY_dot_ij = XY_dot[0]
                a=(n*XY_dot_ij)-(X_sum_ij*Y_sum_ij))/((n*X_pow_2_sum_ij)-(X_sum_ij)*(X_sum_ij))
                b=(Y_sum_ij-a*X_sum_ij)/(n)
                error_matrix[i][j]=abs((Y_pow_2_sum_ij)+(a*a)*X_pow_2_sum_ij)+(n*(b*b))+(2*a*b*X_sum_ij)-(2*a
            else:
                n = j+1
                #calc individual params
                X_sum_ij = X_sum[n-1] - X_sum[i-1]
                Y_sum_ij = Y_sum[n-1] - Y_sum[i-1]
                X_pow_2_sum_ij = X_pow_2_sum[n-1] - X_pow_2_sum[i-1]
                Y_pow_2_sum_ij = Y_pow_2_sum[n-1] - Y_pow_2_sum[i-1]
                XY_dot_ij = XY_dot[n-1] - XY_dot[i-1]
                a=(n*XY_dot_ij)-(X_sum_ij*Y_sum_ij))/((n*X_pow_2_sum_ij)-(X_sum_ij)*(X_sum_ij))
                b=(Y_sum_ij-a*X_sum_ij)/(n)
                error_matrix[i][j]=abs((Y_pow_2_sum_ij)+(a*a)*X_pow_2_sum_ij)+(n*(b*b))+(2*a*b*X_sum_ij)-(2*a
    return error_matrix
```

$\rightarrow O(n)$

$\downarrow O(1)$

$\rightarrow O(n^2)$

Calculate Error Matrix variables

```
In [3]: def preprocess_arrays(n_list,x_list,y_list):
    #let's find cumulative sum for all the points error calculation variables
    #consider points from 0 to n_list-1 , sum[0] means sum till the point 0
    X_sum = np.zeros(n_list)
    Y_sum = np.zeros(n_list)
    X_pow_2_sum = np.zeros(n_list)
    Y_pow_2_sum = np.zeros(n_list)
    XY_dot = np.zeros(n_list)
    #let's calculate the cumulative sum for all these variables in a single for loop
    # let's calculate it for 0
    X_sum[0] = x_list[0]
    Y_sum[0] = y_list[0]
    X_pow_2_sum[0] = x_list[0]*x_list[0]
    Y_pow_2_sum[0] = y_list[0]*y_list[0]
    XY_dot[0] = x_list[0]*y_list[0]
    for i in range(0,n_list):
        X_sum[i] = X_sum[i-1]+x_list[i]
        Y_sum[i] = Y_sum[i-1]+y_list[i]
        X_pow_2_sum[i] = X_pow_2_sum[i-1] + (x_list[i]*x_list[i])
        Y_pow_2_sum[i] = Y_pow_2_sum[i-1] + (y_list[i]*y_list[i])
        XY_dot[i] = XY_dot[i-1] + (x_list[i]*y_list[i])
    return X_sum,Y_sum,X_pow_2_sum,Y_pow_2_sum,XY_dot
```

Total = $O(n)$

$\downarrow O(1)$

$\downarrow O(1)$

$\downarrow O(1)$

$\rightarrow O(n)$

∴ The total time complexity for my algorithm is $O(n^2)$ for a single dataset record