

# Analysis of Mark Sort Algorithm for Database column swaps

```

48 void mrkSort(Array2D *a, int col){
49     for(int i=0; i<a->rows-1; i++){
50         for(int j=i+1; j<a->rows; j++){
51             if(a->data[i][col]>a->data[j][col]){
52                 for(int k=0; k<a->cols; k++){
53                     int temp=a->data[i][k];
54                     a->data[i][k]=a->data[j][k];
55                     a->data[j][k]=temp;
56                 }
57             }
58         }
59     }
60 }
    
```

Equation representing the code ( $O \rightarrow$  Operations)

$O_0 \rightarrow$  Operations outer loop

$O_i \rightarrow$  " inner loop

$O_{if} \rightarrow$  " after if

$O_s \rightarrow$  " swap

$R = \text{Row}, C = \text{col}, L, J = \text{indexes}$

then

$$O + \sum_{i=0}^{R-2} \left( O_0 + \sum_{j=i+1}^{R-1} \left( O_i + P \left( O_{if} + \sum_{k=0}^{C-1} O_s \right) \right) \right)$$

Let  $O_i + P(O_{if} + CO_s) = S$

Substituting

$$O + \sum_{i=0}^{R-2} \left( O_0 + \sum_{j=i+1}^{R-1} S \right)$$

Next round of simplifications with goal to obtain a function dependent on  $R$

$$f(R) = C_0 + C_1 R + C_2 R^2$$

to continue

$$f(R) = 0 + \sum_{i=0}^{R-2} \left( \underbrace{0_0 + \sum_{j=i+1}^{R-1} s}_{0_0 + ((R-1) - (i+1) + 1)s} \right)$$
$$0_0 + (R-1)s - is$$

for the next loop

$$f(R) = 0 + \sum_{i=0}^{(R-2)} (0_0 + (R-1)s - is)$$
$$= 0 + (R-1)0_0 + (R-1)s - \frac{(R-2)(R-1)}{2}s$$

which simplifies to

$$f(R) = (0 - 0_0) + (0_0 - s/2)R + (s/2)R^2$$

Therefore

$$C_0 = 0 - 0_0$$

$$C_1 = 0 - s/2 \quad \text{where } s = 0_i + P(0_{if} + C_0 s)$$

$$C_2 = s/2$$

finally

$$f(R) = C_0 + C_1 R + C_2 R^2$$

$\frac{3}{8}$   
Q.E.D.

What does this say about how this algorithm scales, i.e.  $R \rightarrow \infty$

$$f(R) = O(g(R))$$

$$= \left\{ f(R) : R, R_0 \in \mathbb{Z}, C \in \mathbb{R}; \right. \\ \left. 0 \leq f(R) \leq Cg(R) \forall R \geq R_0 \right\}$$

$$\text{let } g(R) = R^2$$

$$\text{then } 0 \leq C_0 + C_1 R + C_2 R^2 \leq C R^2$$

$$0 \leq \frac{C_0 + C_1 R + C_2 R^2}{R^2} \leq C$$

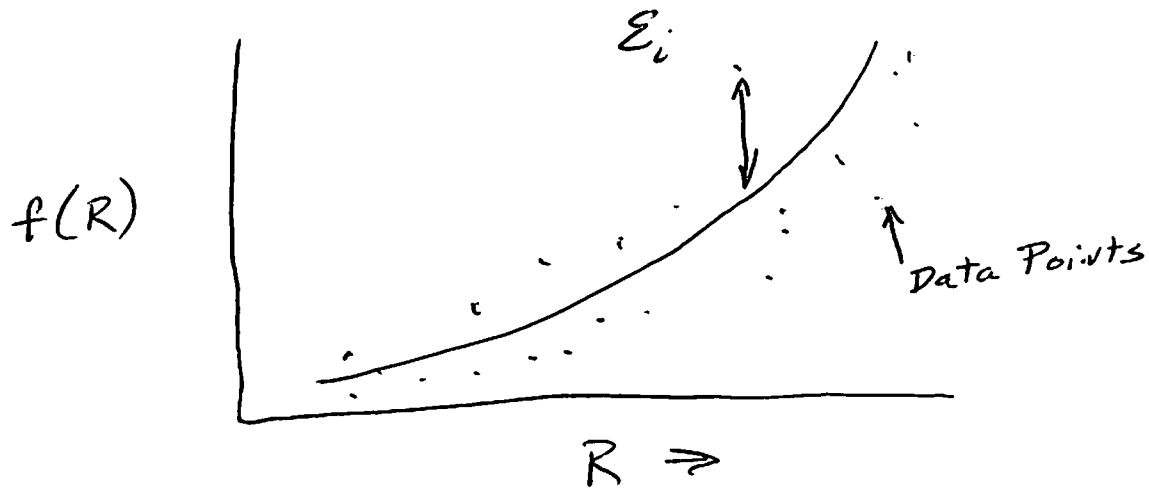
$$0 \leq \frac{C_0}{R^2} + \frac{C_1}{R} + C_2 \leq C$$

$$\lim_{R \rightarrow \infty} 0 \leq \cancel{\frac{C_0}{R^2}}^0 + \cancel{\frac{C_1}{R}}^0 + C_2 \leq C$$

$$\text{then } f(R) = O(R^2)$$

Q.E.D.

Derivation of Least Squares Curve Fit for empirical data derived from timing or operational analysis.



Let

$m \rightarrow$  Number of data points  
 $N \rightarrow$  Number of coefficients to solve for from empirical data

$$m > N$$

$R_{m \times N} \rightarrow$  Matrix relating size or number of rows run vs. the function to curve fit.

$\underline{f}_{m \times 1} \rightarrow$  Vector containing the empirical number of operations or time algorithm took to run.

$\underline{\epsilon}_{m \times 1} \rightarrow$  Vector representing error between curve to fit and actual data points

$\underline{C}_{N \times 1} \rightarrow$  Proposed Vector of coefficients to find relating to the type of function.

Equation which defines the curve you're trying to fit with its associated errors.

$$\underline{f}_{M \times 1} = \underline{R}_{M \times N} \underline{C}_{N \times 1} + \underline{\epsilon}_{M \times 1}$$

Goal: Minimize the errors to find the best curve fit.

Rearranging in terms of the errors

$$\underline{\epsilon} = \underline{f} - \underline{R} \underline{C}$$

Square the error which becomes our cost function

$$J(\underline{C}) \Rightarrow \underline{\epsilon}_{1 \times M}^T \underline{\epsilon}_{M \times 1} = (\underline{f} - \underline{R} \underline{C})^T (\underline{f} - \underline{R} \underline{C})$$

$T \rightarrow$  represents the transpose

$$= \underline{f}^T \underline{f} - \underline{f}^T \underline{R} \underline{C} - \underline{C}^T \underline{R}^T \underline{f} + \underline{C}^T \underline{R}^T \underline{R} \underline{C}$$

$$= \underline{f}^T \underline{f} - 2 \underline{f}^T \underline{R} \underline{C} + \underline{C}^T \underline{R}^T \underline{R} \underline{C}$$

Note:

$$J(\underline{c}) = \underline{f}^T \underline{f} - 2 \underline{f}^T \underline{R} \underline{c} + \underline{c}^T \underline{R}^T \underline{R} \underline{c}$$

is the same as

$$= \underline{f}^T \underline{f} - 2 \underline{c}^T \underline{R}^T \underline{f} + \underline{c}^T \underline{R}^T \underline{R} \underline{c}$$

We would like to minimize the cost function by taking the derivative and setting equal to the zero vector

$$\underline{\phi} = \frac{\partial J(\underline{c})}{\partial \underline{c}} = \frac{\partial (\underline{f}^T \underline{f} - 2 \underline{c}^T \underline{R}^T \underline{f} + \underline{c}^T \underline{R}^T \underline{R} \underline{c})}{\partial \underline{c}}$$

$$\underline{\phi}_{N \times 1} = - \underbrace{2 \underbrace{\underline{R}^T}_{N \times M} \underbrace{\underline{f}}_{M \times 1}}_{N \times 1} + \underbrace{2 \underbrace{\underline{R}^T \underline{R}}_{N \times M \times M \times N}}_{N \times 1} \underline{c}_{N \times 1}$$

Which gives

$$\underline{R}^T \underline{R} \underline{c} = \underline{R}^T \underline{f}$$

and finally

$$\underline{c} = (\underline{R}^T \underline{R})^{-1} \underline{R}^T \underline{f}$$

Q.E.D.