# Mapreduce, gn, cp, pagerank, fm, dgim Numericals

Big Data Analytics (University of Mumbai)

Q. Mapreduce.

$$A = \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix} \quad 3\times3 \qquad B = \begin{bmatrix} 10 \\ 11 \\ 12 \end{bmatrix} \quad 3\times1$$

$$A_{ij} = \begin{bmatrix} 3 & 4 \\ 7 & 2 \\ 5 & 9 \end{bmatrix} \quad 3\times2 \qquad B_{jk} = \begin{bmatrix} 3 & 1 & 5 \\ 6 & 9 & 7 \end{bmatrix} \quad 2\times3.$$

$i = 3, \quad j = 2, \quad k = 3.$

① Map.    (Matrix A)       (Matrix B).

| $(i, k)$ | $(A, j, A_{ij})$ | $(i', k)$ | $(B, k, B_{jk})$ | |
|---|---|---|---|---|
| (1,1) | (A, 1, 3) | (1,1) | (B, 1,  ) | 3 |
| (1,2) | (A, 1, 3) | (1,2) | (B, 1,  ) | 1 |
| (1,3) | (A, 1, 3) | (1,3) | (B, 1,  ) | 5 |
| (1,1)(2,1) | (A, 2, 4) | (1,1) | (B, 2,  ) | 6 |
| (1,2)(2,2) | (A, 2, 4) | (1,2) | (B, 2,  ) | 9 |
| (1,3)(2,3) | (A, 2, 4) | (1,3) | (B, 2,  ) | 7 |
| (2,1)(3,1) | (A, 1, 7) | (2,1) | (B, 3,  ) | 3 |
| (2,2)(3,2) | (A, 7, 7) | (2,2) | (B, 3,  ) | 1 |
| (2,3)(3,3). | (A, 1, 7) | (2,3) | (B, 3,  ) | 5 |
| (2,1) | (A, 2, 2) | (2,1) | (B, 7, 6) | 6 |
| (2,2) | (A, 2, 2) | (2,2) | (B, 7, 6) | 9 |
| (3,3) | (A, 2, 2) | (2,3) | (B, 7, 6) | 7 |
| (3,1) | (A, 1, 5) | (3,1) | (B, 2, 9) | 3 |
| (3,2) | (A, 1, 5) | (3,2) | (B, 7, 9) | 1 |
| (3,3) | (A, 1, 5) | (3,3) | (B, 2, 9) | 5 |
| (3,1) | (A, 2, 9) | (3,1) | (B, 3,  ) | 6 |
| (3,2) | (A, 2, 9) | (3,2) | (B, 3,  ) | 9 |
| (3,3) | (A, 2, 9) | (3,3) | (B, 3,  ) | 7 |

**Reducer.**

$$(X, j, x_{ij}/x_{jk})$$

**Shuffle**

| $(i,k)$ | | |
|---|---|---|
| $(1,1)$ | $(A,1,3)$, | $(A,2,4)$ |
| | $(B,1,3)$ | $(B,2,1)$ |
| $(1,2)$ | $(A,1,3)$ | $(A,2,4)$ |
| $(1,2)$ | $(B,1,1)$ | $(B,2,9)$ |
| $(1,3)$ | $(A,1,3)$ | $(A,2,4)$ |
| | $(B,1,5)$ | $(B,2,7)$ |
| $(2,1)$ | $(A,1,7)$ | $(A,2,2)$ |
| | $(B,3,3)$ | $(B,1,6)$ |
| $(2,2)$ | $(A,1,7)$ | $(A,2,2)$ |
| | $(B,3,1)$ | $(B,1,9)$ |
| $(2,3)$ | $(A,1,7)$ | $(A,2,2)$ |
| | $(B,3,5)$ | $(B,1,7)$ |
| $(3,1)$ | $(A,1,5)$ | $(A,2,9)$ |
| $(3,1)$ | $(B,2,3)$ | $(B,3,6)$ |
| $(3,2)$ | $(A,1,5)$ | $(A,2,9)$ |
| | $(B,2,1)$ | $(B,3,9)$ |
| $(3,3)$ | $(A,1,5)$ | $(A,2,9)$ |
| | $(B,2,5)$ | $(B,3,7)$ |

**Reducer.**

$I_{ij,k}$

② Reducer — for matching $(i,k)$ keep diff x, vertically multiply the values & add the product horizontally.

| | | |
|---|---|---|
| $(1,1)$ | $(3\times3 + 4\times1)$ | $= 13$ |
| $(1,2)$ | $(3\times1 + 4\times9)$ | $= 39$ |
| $(1,3)$ | $(3\times5 + 4\times7)$ | $= 43$ |
| $(2,1)$ | $(7\times3 + 2\times6)$ | $= 33$ |
| $(2,2)$ | $(7\times1 + 2\times9)$ | $= 25$ |
| $(2,3)$ | $(7\times5 + 2\times7)$ | $= 49$ |

$(3,1)$  $(5\times3 + 9\times6)$ $= 87$

$(3,2)$  $(5\times1 + 9\times9)$ $= 86$

$(3,3)$  $(5\times5 + 9\times7)$ $= 88$

$$Product = \begin{array}{c} \\ 1 \\ 2 \\ 3 \end{array} \overset{\displaystyle 1 \quad 2 \quad 3}{\begin{bmatrix} 13 & 39 & 43 \\ 33 & 25 & 49 \\ 87 & 85 & 88 \end{bmatrix}} \qquad \begin{bmatrix} 13 & 33 & 87 \\ 39 & 25 & 86 \\ 43 & 49 & 88 \end{bmatrix}$$

$\overline{\phantom{xxxxx}}\times\overline{\phantom{xxxxx}}$

## * TWO STEP

$$M = \begin{array}{c} 1 \\ 2 \end{array}\overset{\displaystyle 1 \quad 2}{\begin{bmatrix} 2 & 3 \\ 4 & 9 \end{bmatrix}}_{i\times j} \qquad N = \begin{array}{c} 1 \\ 2 \end{array}\overset{\displaystyle 1 \quad 2 \quad 3}{\begin{bmatrix} 5 & 6 & 8 \\ 2 & 4 & 7 \end{bmatrix}}_{j\times k}$$

$i = 2, \quad j = 2, \quad k = 3.$

Step 1  — Mapper.   $\left\{ \begin{array}{l} j\,(M, \ell, M_{ij}) \\ j\,(N, u, N_{jk}) \end{array} \right\}$

| $j\,(M, i, M_{ij})$ | $j\,(N, k, N_{jk})$ |
|---|---|
| 1 (M, 1, 2) | 1 (N, 1, 5) |
| 1 (M, 2, 4) | 1 (N, 2, 6) |
| 2 (M, 1, 3) | 1 (N, 3, 8) |
| 2 (M, 2, 9) | 2 (N, 1, 2) |
|  | 2 (N, 2, 4) |
|  | 2 (N, 3, 7) |

Reducer   $j\,(i, k, M_{ij} \times N_{jk})$

| | |
|---|---|
| 1 (1, 1, 2×5) $\Rightarrow$ | 1 (1, 1, 10) |
| 1 (2, 1, 4×5) $\Rightarrow$ | 1 (2, 1, 20) |
| 1 (1, 2, 2×6) $\Rightarrow$ | 1 (1, 2, 12) |
| 1 (2, 2, 4×6) $\Rightarrow$ | 1 (2, 2, 24) |
| 1 (1, 3, 2×8) $\Rightarrow$ | 1 (1, 3, 16) |
| 1 (2, 3, 4×8) $\Rightarrow$ | 1 (2, 3, 32) |

2 (1,1 , 3×2 ) ⇒ 2(1,1, 6)
2 (2,2, 9×2 ) ⇒ 2(2,2,18.)
2 (1,2, 3×4 ) ⇒ 2(1,2, 12)
2 (2,2, 9×4 ) ⇒ 2(2,2, 36)
2 (1,2, 3×7 ) ⇒ 2(1,3, 21)
2 (2,3, 9×7) ⇒ 2(1,2, 63)

Step 2 - Mappes.
(j $[(i_1, k_1, v_1), (i_2, k_2, v_2) \cdots$ ])

(1 [(1,1,10), (1,2,12), (1,3,16),
(2,1, 20), (2,2,24), (2,3,32) ] )

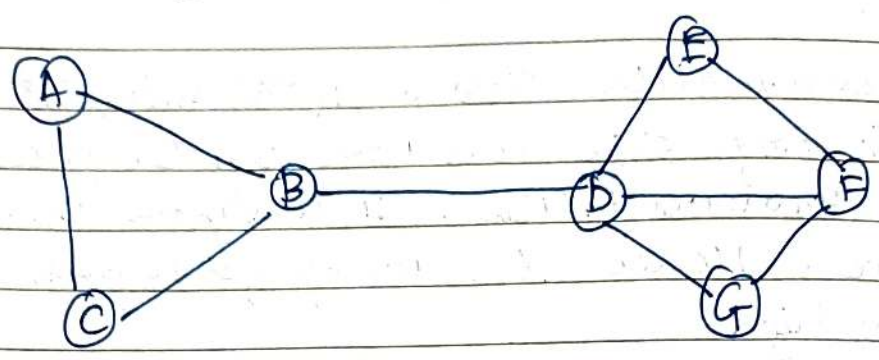(2 [(1,1 6 ), (1,2, 12 ), (1,3,21),
(2,1,18), (2,2,36), (2,3,63) ].

Reducer

For each (i,k), add the corresponding
values.
[[(1,1), 10+6]          [(2,1), 20+18
[(1,2), 12+12]         [(2,2), 24+36
[(1,3), 16+21]]        [(2,3), 32+63]]

P = [16 [16  24  37]
     24 [38  60  95]
     37 ]

B. GN algorithm.

<u>Shortest paths.</u>                                    EB

Path.

AB ⇒ AB, AD, AE, AG, AF, AC ⇒ 45

AC ⇒ AC                              ⇒ 1

BCAD ⇒ AB, AD, AE, AG, AF, AC ⇒ 5

BD ⇒ { BD, BE, BF, BG
        AD, AE, AF, AG        ⇒ 4×3 ⇒ 12
        CD, CE, CF, CG }

DE ⇒ DE, CE
      BE,              } = PG ⇒ 4
      AE

DF ⇒ DF, BF, AF, CF          ⇒ 4.

EF ⇒ { EF, EG^½ } ⇒ ED       ⇒ 1.5

Breaking the edge with highest EB. {BD}.

communities.

$DF \to DE, DF^{1/2}$    $EF \to EF, EG^{1/2}$

$DG \to DG, DF^{1/2}$    $FG \to FG, EG^{1/2}$

EB. i.e,

Algo → ① calculate the total no. of shortest paths from some pt $\alpha$ to some pt $\beta$ that pass through edge E.

② Break the edge with highest EB factor.

③ Repeat till required no. of comm is obtained or all edges in a comm have the same EB factor.

→ Points.

- hierarchical clustering model

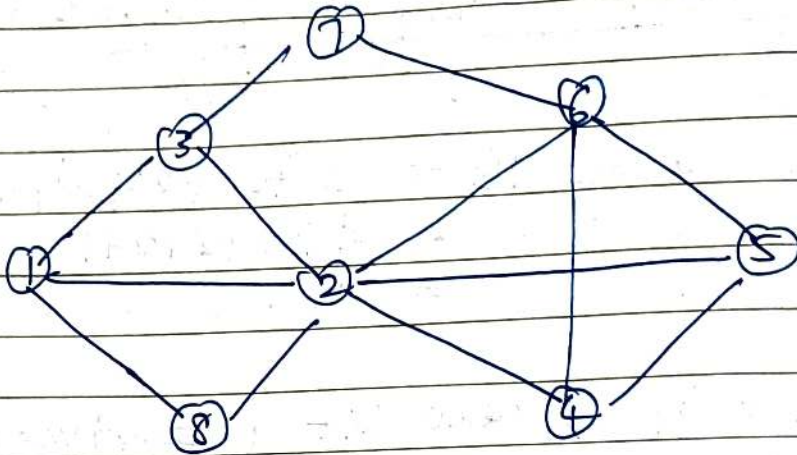- used for ① comm. detection

        ② measure edge betweeness

- $EB_E = \left\{ \begin{array}{l} \text{no of shortest paths from some node} \\ X \text{ & some other node } Y \text{ that passes} \\ \text{through edge E} \end{array} \right.$

- if there are $n$ shortest paths between $X$ & $Y$, then all edges $E_1, E_2 \ldots E_n$ will have $(1/n)$ weight.

- complexity → ① calculation of edge betn $= O(EN)$

             ② algorithm $= O(E^{2N})$

**Q.** Clique percolation method.



Identifying 3 cliques. $\Rightarrow$ 6 3-cliques.

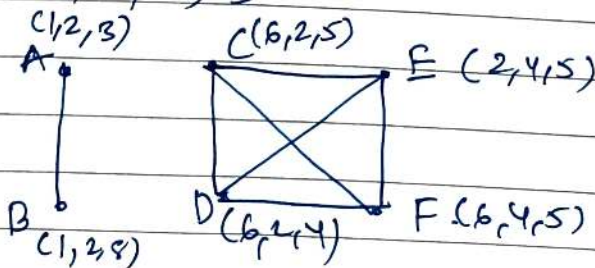A: (1,2,3)          E: (2,4,5)
B: (1,2,8)          F: (6,4,5)
C: (6,2,5)
D: (6,2,4)

Combining cliques into communities.
( cliques should have 2 nodes in common).

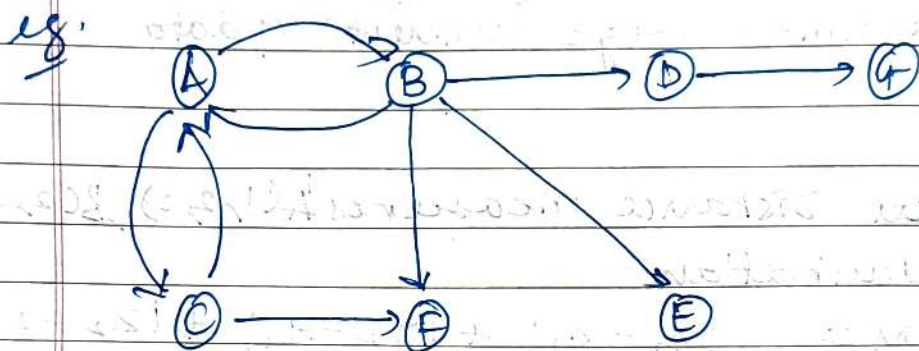A: (1,2,3)     $\Big\}$ $\rightarrow$ $c_1$ : (1,2,3,8)
B: (1,2,8)     $\Big\}$

C: (6,2,5)     $\Big\}$ $\rightarrow$ $c_2$ : (2,4,5,6)
D: (6,2,4)     $\Big\}$

E: (2,4,5)     $\Big\}$ $\rightarrow$ $c_3$ : ()
F: (6,4,5)     $\Big\}$

# Q. Pagerank Algorithm

Purpose is to rank the webpages based on the no of incoming & outgoing links in that page. A more important page will have larger no-of such links.

eg.



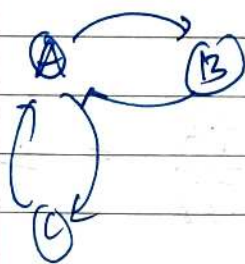In this case, E & G are deadends.

I becomes deadend.

Transition Matrix $M_T =$

|   | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| A | 0 | $\frac{1}{4}$ | $\frac{1}{2}$ | 0 | 0 | 0 | 0 |
| B | $\frac{1}{2}$ | 0 | 0 | 0 | 0 | 0 | 0 |
| C | $\frac{1}{2}$ | 0 | 0 | 0 | 0 | 0 | 0 |
| D | 0 | $\frac{1}{4}$ | 0 | 0 | 0 | 0 | 0 |
| E | 0 | $\frac{1}{4}$ | 0 | 0 | 0 | 0 | 0 |
| F | 0 | $\frac{1}{4}$ | $\frac{1}{2}$ | 0 | 0 | 0 | 0 |
| G | 0 | 0 | 0 | 1 | 0 | 0 | 0 |

now $M_T =$



$$M_T = \begin{array}{c} A \\ B \\ C \end{array} \begin{bmatrix} 0 & 1 & 1 \\ \frac{1}{2} & 0 & 0 \\ \frac{1}{2} & 0 & 0 \end{bmatrix} \quad , \quad v_1 = \begin{bmatrix} \frac{V}{3} \\ \frac{V}{3} \\ \frac{V}{3} \end{bmatrix}$$

initial rank

To handle spider trap( group of pages that have no outgoing links.

we introduce teleport factor into the eqn. ($\beta$=0.85). we assume that the web server will visit a page with $\alpha$ prob. & start fresh browsing with $(1-\beta)$ prob. i.e. teleport to some random page with $(1-\beta)$ prob.

new page rank $\Rightarrow$

$$v' = \beta mv + (1-\beta) \cdot e/n$$

iteration 0.

$$v_1' = 0.85 \begin{bmatrix} 0 & 1 & 1 \\ \frac{1}{2} & 0 & 0 \\ \frac{1}{2} & 0 & 0 \end{bmatrix} \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} + 0.15 \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix}$$

$$v_1' = \frac{1}{120} \begin{bmatrix} 74 \\ 23 \\ 23 \end{bmatrix}$$

iteration 1

$$v_2' = \frac{0.85}{120} \begin{bmatrix} 0 & 1 & 1 \\ \frac{1}{2} & 0 & 0 \\ \frac{1}{2} & 0 & 0 \end{bmatrix} \begin{bmatrix} 74 \\ 23 \\ 23 \end{bmatrix} + 0.15 \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix}$$

$$= \frac{1}{2400} \begin{bmatrix} 902 \\ 749 \\ 749 \end{bmatrix}$$

$PR_A = \dfrac{902}{2400}$ , $PR_B = \dfrac{749}{2400}$ , $PR_C = \dfrac{749}{2400}$.
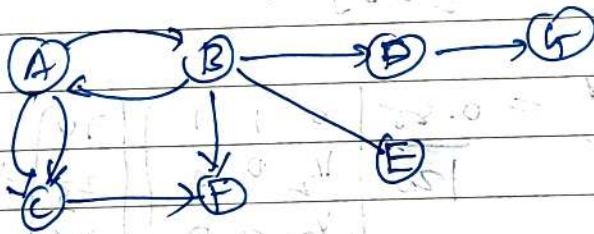
adding deadends. D, E, F, G



$$PR_D = \frac{PR_B}{4} =$$

$$PR_E = \frac{PR_B}{4}$$

$$PR_F = \frac{PR_B}{4}$$

adding G

$$PR_G = PR_D.$$

Q. **FM Algorithm.**

- used to count unique nos in a data stream.
- use hashing ~~them~~ functions
- time complexity is $O(n)$
  memory req - $O(\log(m))$

Algo -
- make an array of stream elements.
- take a hash function $\{h(n) = (5n+1)\%4$
- find hash values of~~or~~ each element
- convert ~~hash val~~ binary of each hash value
- count no of zeros in each bin. value
- ~~mss of~~ find max no of trailing zeros.
- no. of unique values in stream
  $\approx 2^n$ (nearly accurate)
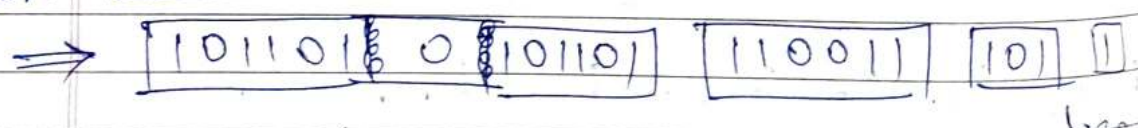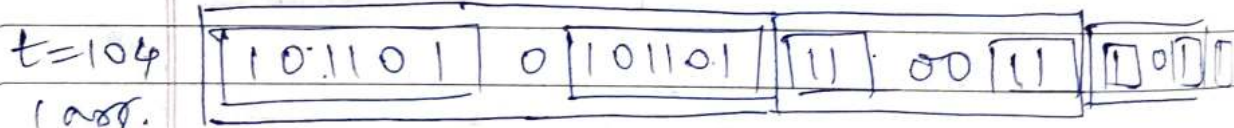
eg. $S = 1, 3, 2, 1, 2, 3, 4, 3, 1, 2, 3, 1$

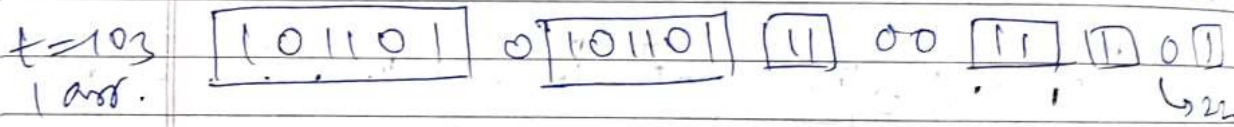| n | 1 | 3 | 2 | 1 | 2 | 3 | 4 | 3 | 1 | 2 | 3 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| h(n) | 2 | 0 | 3 | 2 | 3 | 0 | 1 | 0 | 2 | 3 | 0 | 2 |
| B(n) | 010 | 000 | 010 | 010 | 010 | 000 | 001 | 000 | 010 | 011 | 000 | 010 |
| trailing zeros. | 1 | 2 | 0 | 2 | 0 | 2 | 0 | 2 | 1 | 0 | 2 | 1 |

$r = \max(\text{trailing zeros}) = 2.$

no of unique values $= 2^r = 2^2 = \underline{4} \{1,2,3,4\}$

# DGIM

1 0 1 1 · 0 1　0 1 0 1　1 0 1　1 0 0 1 1

   3     6     9    12   15  18  19.

t=100

1 0 1 1 1 0 1 1 0 1 0 1 0 1 1 1 0 0 1 1

→ 1 0 1 1 0 1 0 1 0 1 1 0 1 1 1 0 0 1 1 ↳19

t=101
1 arr.
1 0 1 1 0 1 0 1 0 1 1 0 1 1 0 0 1 1 1 ↳19

→ 1 0 1 1 0 1 0 1 0 1 1 0 1 1 0 0 1 1 1 ↳20

t=102
0 arr.
1 0 1 1 0 1 0 1 0 1 1 0 1 1 0 0 1 1 1 0 ↳21

t=103
1 arr.
1 0 1 1 0 1 0 1 0 1 1 0 1 1 0 0 1 1 1 0 1 ↳22

t=104
1 arr.
1 0 1 1 0 1 0 1 0 1 1 0 1 1 0 0 1 1 1 0 1 1

→ 1 0 1 1 0 1 0 1 0 1 1 0 1 1 0 0 1 1 1 0 1 1 ↳23

t=105
1 arr
1 0 1 1 0 1 0 1 0 1 1 0 1 1 0 0 1 1 1 0 1 1 1

|←————— WINDOW —————→|

→ 1 0 1 1 0 1 0 1 0 1 1 0 1 1 0 0 1 1 1 0 1 1 1

   23                     22    21  20 20

largest bucket is partly into the window.

calculated

total no of 1's $= \frac{1}{2}(2^3) + 2^2 + 2^1 + 2^0 + 2^0$.

$$= \frac{1}{2}(8) + 4 + 2 + 1 + 1$$
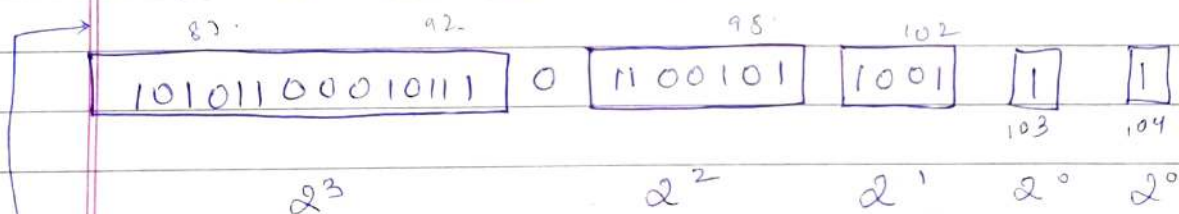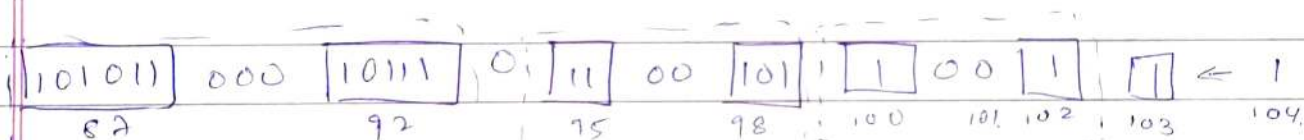
total no of 1's $= 12$

Actual no of 1's $= 16$
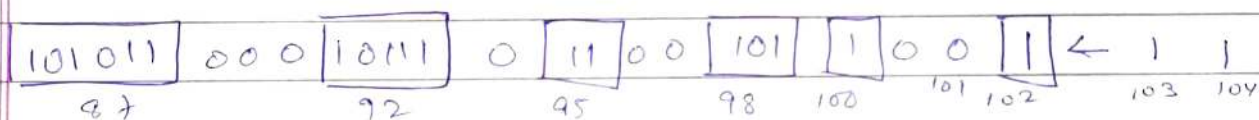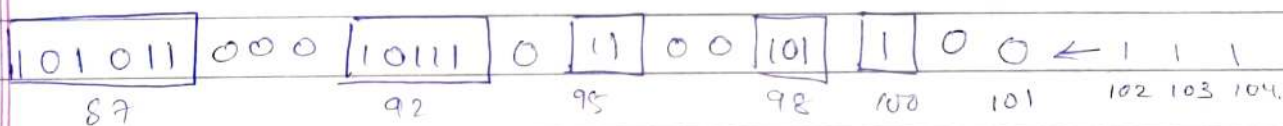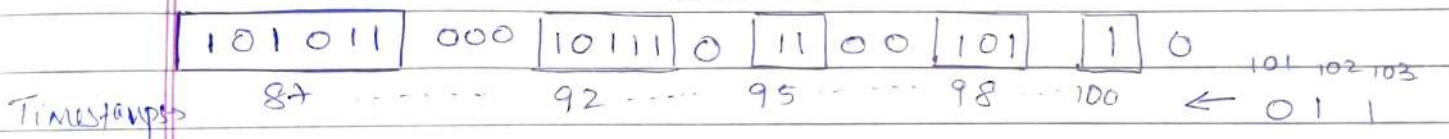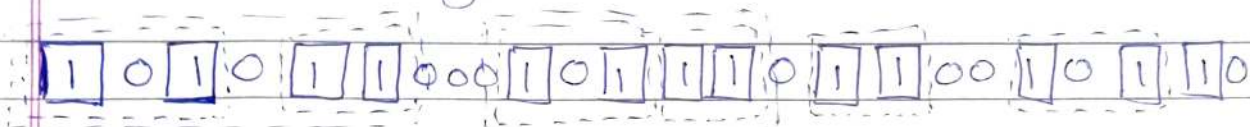
$\therefore$ accuracy $= \frac{\Delta(1's)}{A.V.} \times 100 = \frac{4}{16} \times 100$

$= 25\%$.

# * DGIM

1 0 1 0 1 1 0 0 0 1 0 1 1 1 0 1 1 0 0 1 0 1 1 0

N = 24 (window size).

1 0 1 0 1 1 0 0 1 0 1 1 1 0 1 1 0 0 1 0 1 1 0

| 1 0 1 0 1 1 | 0 0 0 | 1 0 1 1 1 | 0 | 1 1 | 0 0 | 1 0 1 | | 1 | 0 | 0 1 1
| 87 | | 92 | | 95 | | 98 | 100 | | 101 102 103

Timestamps

| 1 0 1 0 1 1 | 0 0 0 | 1 0 1 1 1 | 0 | 1 1 | 0 0 | 1 0 1 | | 1 | 0 0 ← 1 1 1
| 87 | | 92 | | 95 | | 98 | 100 | 101 | 102 103 104.

| 1 0 1 0 1 1 | 0 0 0 | 1 0 1 1 1 | 0 | 1 1 | 0 0 | 1 0 1 | | 1 | 0 0 | 1 ← 1 1
| 87 | | 92 | | 95 | | 98 | 100 | 101 102 | 103 104

| 1 0 1 0 1 1 | 0 0 0 | 1 0 1 1 1 | 0 | 1 1 | 0 0 | 1 0 1 | 1 | 0 0 | 1 | 1 ← 1
| 87 | | 92 | | 95 | | 98 | 100 | 101 102 | 103 | 104.

| 87 | 92 | 95 | 102

| 1 0 1 0 1 1 0 0 0 1 0 1 1 1 | 0 | 1 1 0 0 1 0 1 | 1 0 0 1 | 1 | 1
| | | | | 103 | 104

$2^3$      $2^2$   $2^1$   $2^0$   $2^0$

if current ts - leftmost bucket ts of window < N, continue
eg. 103 - 87 = 16 < 24 ∴ continue.
if greater or equal, stop.

How many 1's are there in the last 20 bits?

$$= \frac{2^3}{2} + 2^2 + 2^1 + 2^0 + 2^0$$

$$= 4 + 4 + 2 + 1 + 1 = \underline{12}$$