

7/12/22

## Query Optimization: କେବଳ ଏକେ କ୍ଷେତ୍ରରେ

କବଳିତ କାହାର କ୍ଷେତ୍ରରେ କ୍ଷେତ୍ରର କବଳିତ କ୍ଷେତ୍ରର  
ଫିଲ୍ଡର୍ କାହାର - every row କେବଳ operation

⇒ ଏକାକି କାମରେ, ଅନ୍ତର୍ଭବାରେ, କ୍ଷେତ୍ରରେ

disk read - write, process

- ଯଦି କ୍ଷେତ୍ରର ଆଇଶାର୍କ କବଳିତ କ୍ଷେତ୍ରର  
କବଳିତ କ୍ଷେତ୍ରର, ତାହାର କବଳିତ କ୍ଷେତ୍ରର  
କବଳିତରେ।
- Linear Search କବଳିତ  $O(n)$  କିମ୍ବା block  
କ୍ଷେତ୍ରରେ, Binary search କବଳିତ  $\log(n)$ .
- ଯଦି primary key କିମ୍ବା କବଳିତ କ୍ଷେତ୍ରର  
column କବଳିତ କବଳିତ, ତଥାତ

P.T.O.

$O(\log n)$  time merge

• Primary index ওঁ clustering property

যাই, Secondary index প্রতি লাই, দ্বারা কোন  
PI এবং থেকে সুজ্ঞা

→ Primary index ফর্ম প্রিমারি key র,  
কাউন্ট ক'রি block কৃত্যতে কো?

$n_r$  = number of tuple

$b_r$  = no. of blocks

$f_r$  = elements in a single block

• جلیل علی خان سعید پور احمد فتحی، سید علی احمدی، سید علی احمدی

جعفر بن محبث

1910S

RFID 000,001 3210. " T :  
RFID 000,001 3210. " T :

$V(A, R) = R$  relation on  $A$  strong

Chapel Hill

$$Sc(A, n) = \sqrt{\frac{\sigma_n}{V(A, n)}} ; \quad b_n = \left\lceil \frac{n}{f_n} \right\rceil$$

### Cardinality

$$24 \text{ ft search log (bn)} + \left[ \frac{\text{Sec}}{\text{ft bn}} \right]$$

pointer is always at first point

$$\therefore \text{ans} = 10 + \log_2\left(\frac{10000}{20}\right) = 10 + 9 = 19 \text{ different blocks}$$

$$\frac{10000}{20} = 500$$

$$\frac{500}{200} = 2.5$$

$$\log_2(2.5) = \frac{10000}{500} = 20$$

Sol:

Select bname = "PUP" from ac;

$$mac = 10,000$$

$$\checkmark (b.name, ac) = 50 \quad (\text{unique}, 50 \geq 1)$$

$$\checkmark (bal, ac) = 500 \quad (\text{unique}, 500 \geq 1)$$

$$\checkmark (b.name, ac) = 500 \quad (\text{unique}, 500 \geq 1)$$

total bname are  
500

$$\text{math} : \frac{foc}{ac} = 20$$

• associativity on key (?) same primary key.

# B+ tree DBP,

$$\text{Fanout, } F = 20$$

$$\log_2 \left( \frac{50,000}{20} \right) + \frac{20}{20} + 1$$

11

2

+ 1

10

↓  
DBP  
shows  
equivalent  
key

Sc  
R

$$\frac{10,000/500}{20}$$

(DBP DB)

DB 2000 logical 2000, 1000 physical for DB

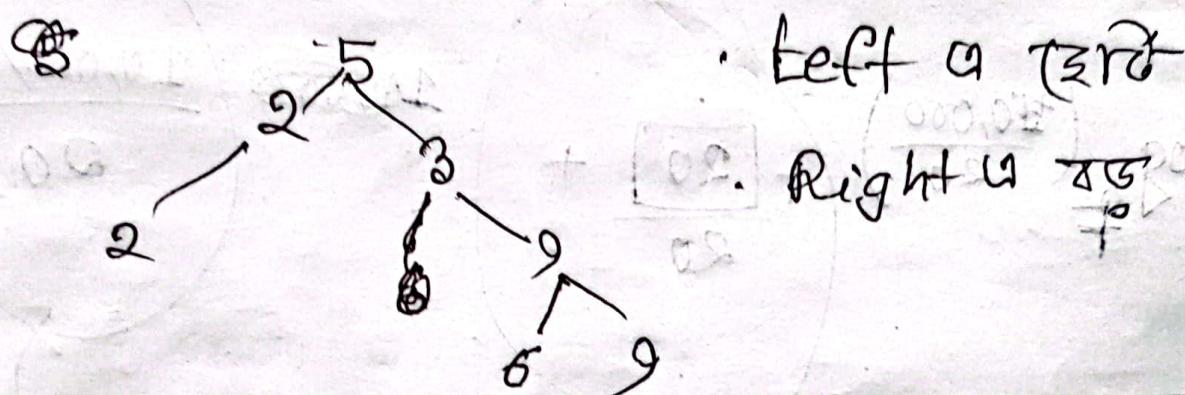
No need to add extra 1000 to be consistent with DB

Background: Indexing এর বাবে

## • K-dimensional tree ⇒

Database এর জন্য Indexing algorithm কোর্ট  
দরবারি, পদক্ষেপ: B+ tree.

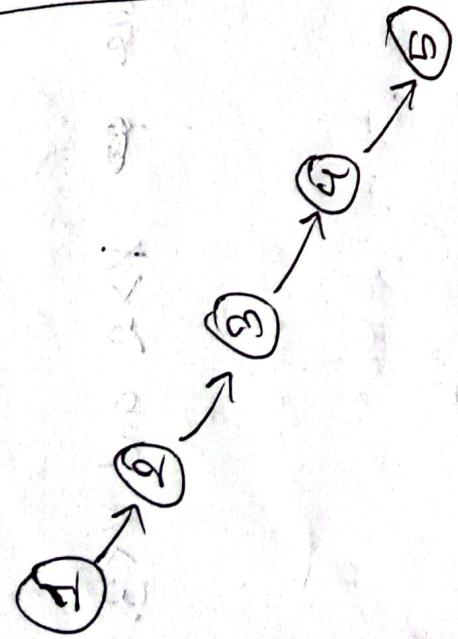
→ Binary tree: 5, 2, 3, 2, 9, 6, 9



এবং অসম্ভব কথা:

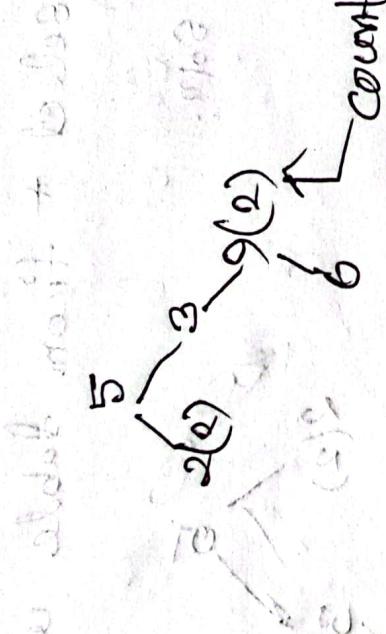
- 1) Repeated data এর ক্ষেত্রে handle করতে  
যাবে না। (তবে যে অসম্ভব সম্মিলন  
counter দ্বারা করা যাবে)
- 2) এটি balanced না হলে worst case  
it gets reduced to linear search.

Insert : 1, 2, 3, 4, 5

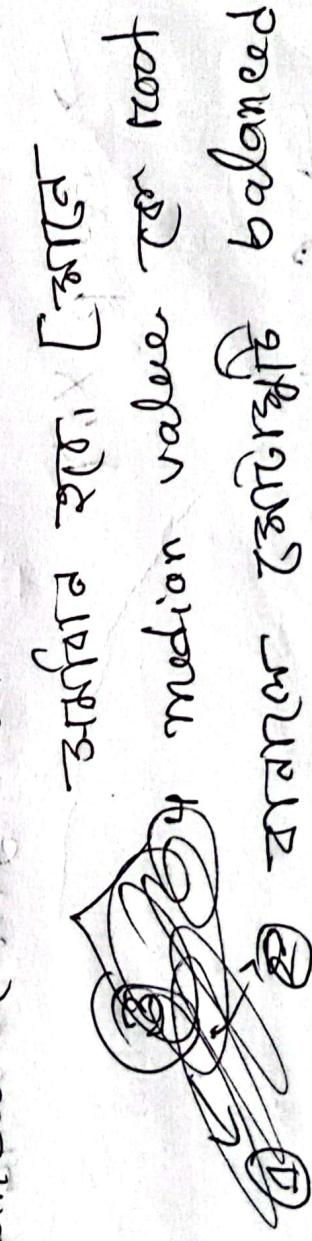


Count Using Counter:

5, 2, 3, 2, 9, 6, 9



Count, year range 75, 214 root value  
or for balanced tree  
balanced tree  
successive numbers  
balanced tree



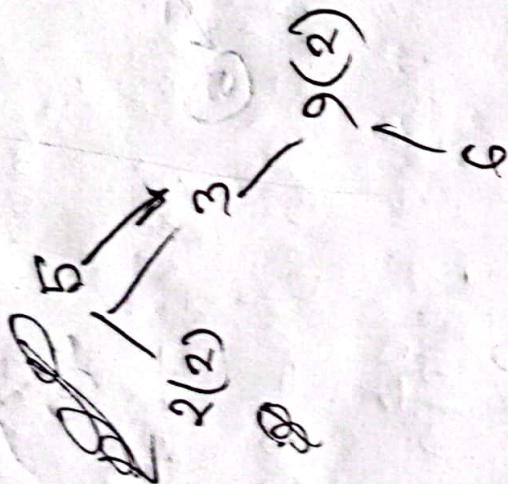
tree insert [1, 2, 4, 4, 10]  
median  
smaller bigger

Range

# Range query in tree:

Select \* from table where  $x > 2$  and  $x < 3$ .

Soln:



# Range query with multiple values,  $(x, y)$ ,  
→ swap  $x, y$  after column

→ swap  $x, y$

→ Grid

13/12/21

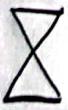
④ T(webpage) = webpage Table

↙ (occurrences, URL) = ~~URL~~ table of ~~occurrences~~  
occurs occur URL

→ size  $10^6 \cdot 2^31$

④ index = join:

i	j	k
1	2	3
2	1	4



④ Index - Lookup: index ~~of~~ BT tree  
natural, ~~not~~ index for search first 2000  
first index for?

first index for?

④ Table Join:

⊗ : natural join

• Left outer join

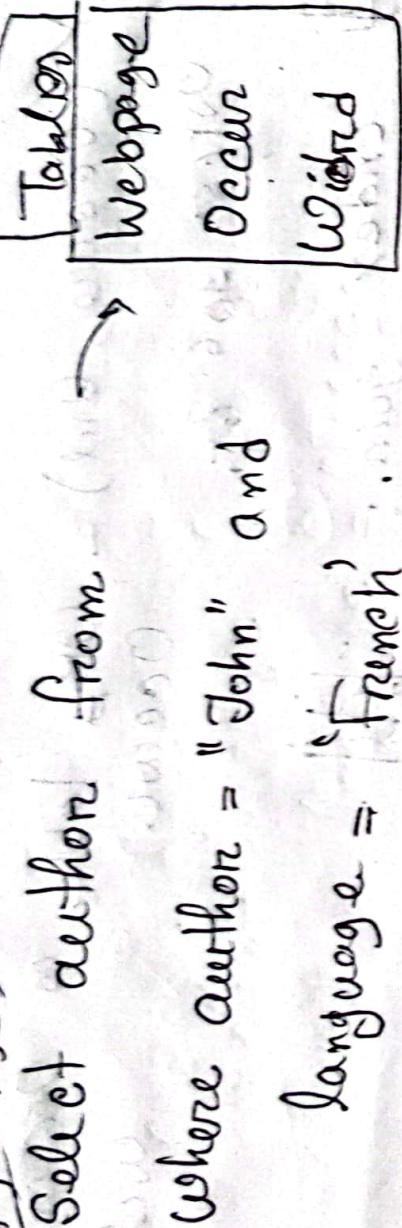
× : Cartesian

⊗ : Left outer join

⊗ : Right outer join

• index - join  $\frac{\text{Table}}{\text{author}}$  = A Table  $\frac{\text{author}}$  word  
 $\text{url} = \text{url}$

Query: page 9



What query miss?

Soln: author and language,

• Select author from webpage where author = "John"

Webpage	Occur	Word
author	• url	• language
index	• index	• index
url	• url	• url

→ Ans: =

$$\frac{10^9}{107} = 100 \text{ or } \text{tuple}$$

- Select Author from Variant editor  
Language is "The lot"

Sol:

100

- Selected author from editor, Author where Author = "John" and webpage.url = author.url.
- Selected author from editor, Author where Author = "John" and webpage.url = author.url.

John John X  $\frac{10^9}{10^{10}}$

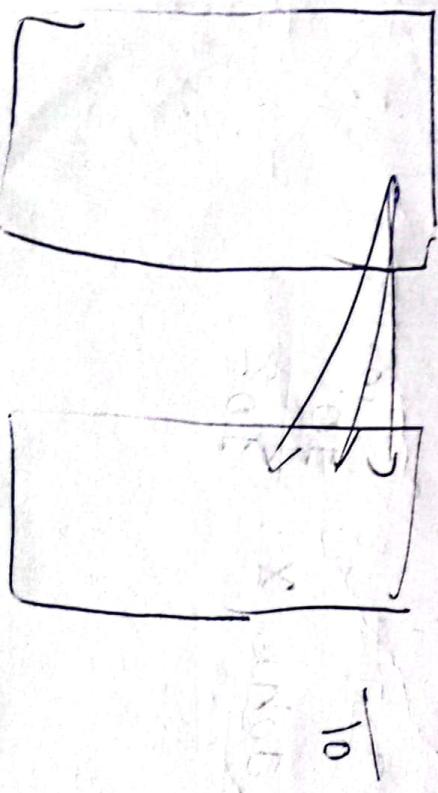
• main - memory - hash - join : Main memory to get block arr, time block access time  $3 \text{ ms}$ .

Select language from word  
culture language = "French"

Ans:  $\frac{10^6}{100} = 10^4$

Final:

$$100 \times \frac{10^3 \times 10^4}{100}$$



$$10^0 + 10^2 \cdot 10^3 + 10^4 + 0 = \text{ans.}$$

$\frac{P(A|B)}{P(B)}$

Likelihood

$P(B|A) P(A)$  ← Prior

Probability

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Posterior

Normalization constant

Posterior probability:

$$\text{Posterior probability} = C \cdot \text{Likelihood}$$

Normalizing factor -  $C$

ज्ञानवाला वर्तमान ज्ञान के साथ प्राप्ति की जानकारी का अनुपात उत्तर देता है। यह अनुपात अधिक लिकेली होती है।

प्राप्तिकरण / Likelihood

$$P(\text{Event}|e) = \frac{P(e|\text{Event}) P(\text{Event})}{P(e)}$$

Posterior probability

प्राप्ति / future.

$$P(\text{Event}) = \text{प्राप्ति का अनुपात } P(\text{Event}) \text{ के लिए } P(\text{Event})$$

தகவல் எல், கால்பனை, prior probability  
குறைபாடு  $P(A|C)$  கூற விரும்புகிறோம்

• enough data என்ற நோக்கம்

Learning algorithm முடிச்சு என்ற நோக்கம்

bayesian எடுத்து கொள்ளுதல்,

$\rightarrow P(C), P(C)$  அதாவத் தொழில் நிலை,

$\rightarrow$  Bias in data:

MLE என்பதை கிடை செய்து  
கொண்ட அறிஞர் MLE க்கால்தான் என்ற நீண்ட பார்வை,  
 $(A) P(A) = P(B)$ .

## Skin detection:

$$T < \frac{P(c|s)}{P(c|ns)} \approx \frac{P(sle)}{P(nsle)}$$

0.4

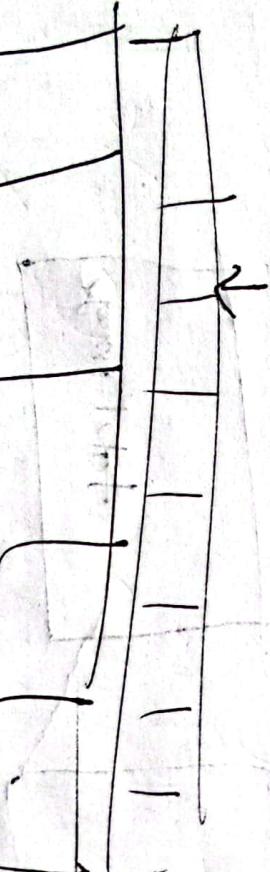
$P(\text{skin area} | \text{bright color}) > T$

$P(\text{no skin} | \text{color})$

255<sup>3</sup> possible colors

Process:

Read image (RGB, 0 - 255)



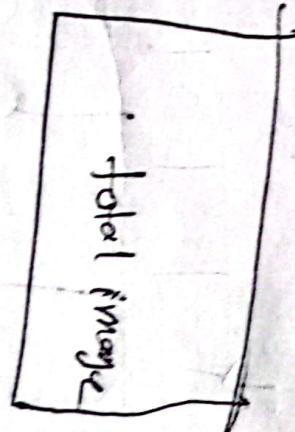
Red

Green

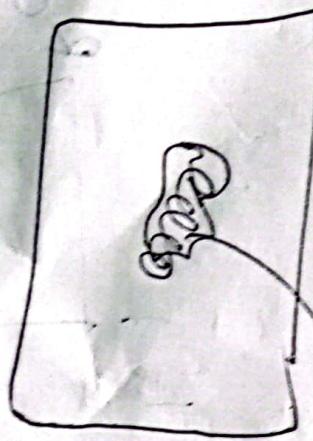
Blue

skins

visual



skin



equivalent file



→ mark feature points

→ mark feature points

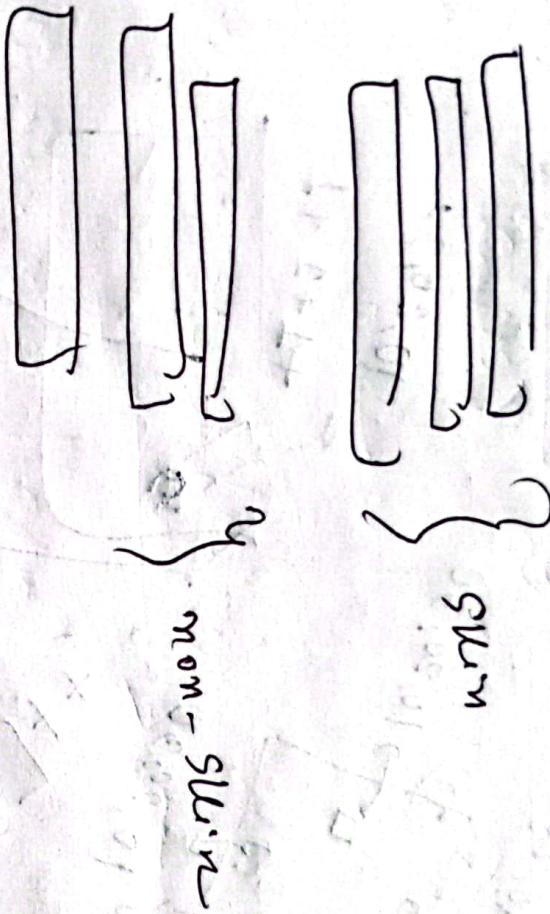
mon-sk.in

→ feature points, convert skin, convert

Ctrl  
186

mon-sk.in

Strong double, 20-28s Tame pair



Test : image after  $\rightarrow$  Colon after  $\rightarrow$  Colon after

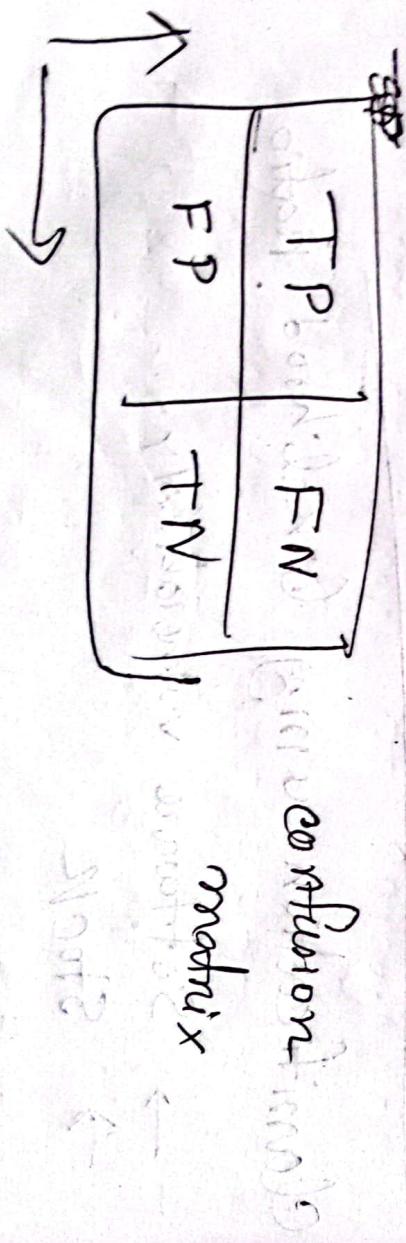
Colon after  $\rightarrow$  Colon after  $\rightarrow$  Colon after

20/12/22

$$\frac{P(c|s)}{P(c|ns)} = \text{Likelihood estimator}$$

### Cross-validation:

- 10 fold cv to reduce error of test
- Randomly 9 crop filter train  
1 crop filter test
- Cross 10 times



~~#10~~  $\rightarrow$  accuracy  $\rightarrow$  avg. accuracy

$\rightarrow$  Confusion matrix  $\rightarrow$  avg. confusion matrix

Task : ML machine learning

repository  $\rightarrow$  classification

$\hookrightarrow$  2 class problem  $\rightarrow$  dataset  $\rightarrow$  jet rate

$\rightarrow$  calculate Accuracy, confusion

matrix trace,

confusion matrix

↓  
few more negative numbers  $\rightarrow$  N

24/12/22

## H.A.: Priori Algorithm:

- Life concern profit maximizing  
and data correlation.

[ DataMining Book 2nd edition.pdf  
Page 236/265 ]

↳ All possible combination of  
frequency to sets, & combination  
frequency to set, therefore  
(parallel) parallel universes

But complexity is very high

Q.

त्रिकोणीय शैम निम्नलिखित ब्रेक्ट

Total algorithm Complexity =  $5! = 120$

पर्याप्त 100 द्वारा अवश्यक, अर्थात् 100! - ग्रन्थे

प्राप्त हुए,

Possible Solution

abc  
bac  
cab

a b c
a b ↙
a c
a b ↙

unique

$O(n!)$

a
b
c
a c
1

b c

a c

a b

b c

..

## Book Solution

- Support Count.
- Satisfies the customer & not about picking the absolute max profit.

↳ Association prior knowledge.

confidence,  $A \Rightarrow B$  अर्थ, A किए

दृष्टि B (2) होता, तो सम्भवता रहती

$$P(A \Rightarrow B) = \frac{sc(A \cup B)}{sc(A)}$$

•  $\frac{\text{संज्ञानात् दिखाया}}{\text{संज्ञानात् दिखाया}}$

$$\therefore A \text{ दृष्टि } B \text{ दृष्टि } = \frac{\text{प्रबलता } A, B \text{ दृष्टि}}{\text{श्रृंखला } A \text{ दृष्टि}}$$

→ Confidence vs Probability

-Likelihood

$$\sum p = 1$$

$$\sum c \neq 1$$

→ Likelihood vs Probability

Likelihood → often called似然性  
Probability → 概率 Distribution.

A

$I_1 I_2 \Rightarrow I_B$

$$= \frac{sc(I_1 I_2) \wedge I_B}{sc(I_1 I_2)}$$

Each  $I_i$  ( $i \in A$ ) is

$I_1 I_2^*, I_B^*, I_1 I_2 I_B^*$ .

$I_1, I_2$

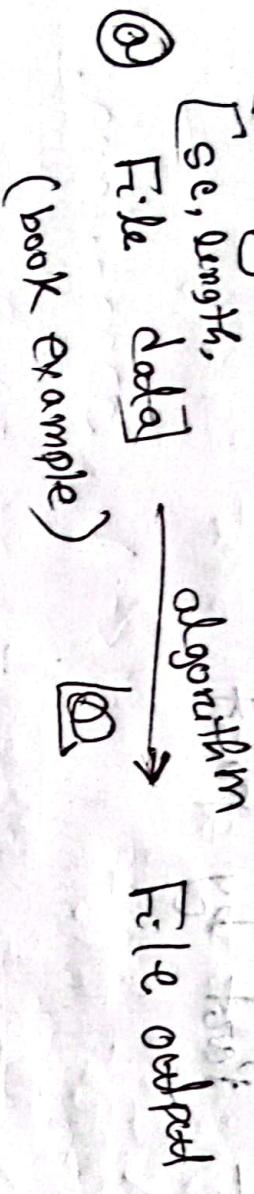
$$= sc(I_1)$$

$$T_1 \wedge T_5 \Rightarrow T_2 = \frac{(T_1 \wedge T_5 \wedge T_2)_{SC}}{(T_1 \wedge T_5)_{SC}}$$

$$= \frac{2}{2} = 100\%$$

∴ We have 100% confidence on  $T_1 \wedge T_5 \Rightarrow T_2$

Coding approach:



③ generate N + criteria query to T2

confidence % calculate matrix result

↳ consider

## Recovery System

We try to make our system as volatile as possible.

- having multiple copy in diff geographical location
- Part by part

But

Defeat start delay or, Defeat consistency in database.

 ACID:

ACID = Atomicity, Consistency, Isolation, Durability

Atomicity: दुनांत एकेजे transaction युद्धाते  
रहूँ, नाही उल्लंघन का,

A gives B 50. } both रद्दमाय,  
B gets 50 from A }

Consistency: फिनांसिअल आव्हानात एक विचार मिळाले  
रहै (accounting एक बजेच)

Isolation: 3-4 चे transaction, parallelly  
चालू झाले तरी प्रत्येके नियमांचे अंतर  
serially चालावून अवैध output नाही,

"A ; B के 100 टोकार नियर असू वा के  
A असू 10% नियर" - येथे ~~असू~~ parallelly  
काढू असून नाही,

Durability: एकीफे त्रिष्ठा आणलेले database  
नाही वापरात,

- Roles in idempotent operations

log manager 950 फे दूर,

↳ one Commit, one redo log

Immediate vs Deferred

↳ Immediate dataset change

↳ Transaction manager

After log

UND

Start unit test and start reloop after 1

↳ Recovery: Log & checkpoint

↳ start transaction and

↳ commit transaction

undo - list

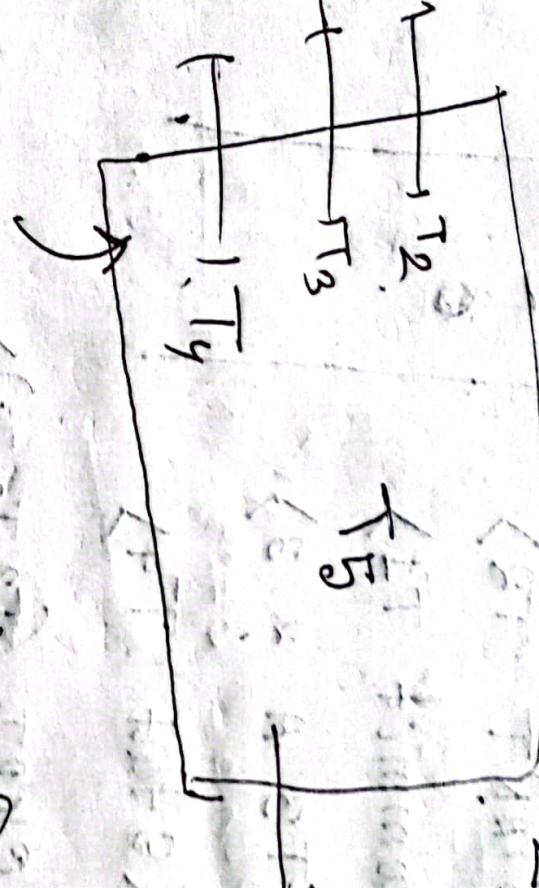
- redo - list

T<sub>3</sub>

T<sub>4</sub>

T<sub>5</sub>

T<sub>6</sub>



ক্ষেত্র পরিবর্তন করে নেওয়া হল।

১০৮৫

→ Trace back

`<END CKP>`

`<T6 D 2000 3000>`

`<START T6>`

`<Commit T1>`

`<T4 C 100 200>`

`<T5 D 1000 2000>`

`<START T5>`

`<Commit T2>`

`<T2 B 10 20>`

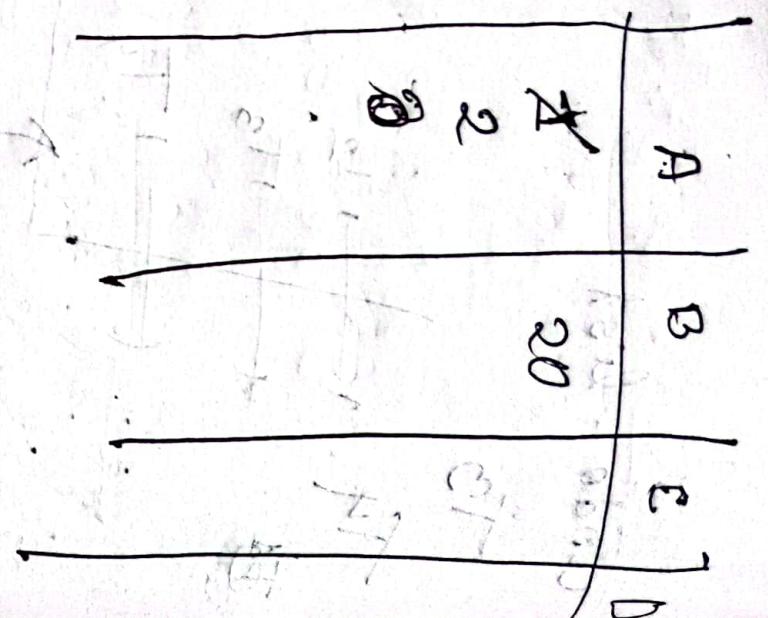
`<CKPT (T2, T3, T4)>`

`<Start T4>`

`<UND <T3 A 2, 3>>`

`<START T2>`  
~~<RED~~  
`<Commit T1>`

`<T1 A 1, 2>`



## Query Optimization Question Solve

a)  $\sigma_{B=25}(R)$

$$Sc(R, B) = \frac{n_R}{V(R, B)} = \frac{100}{10} = 10$$

∴ Number of tuples = 10.

b)  $\sigma_{B=25}$  and  $(C=0)$  ( $R$ )

$$Sc(R, B) = \frac{n_R}{V(R, B)}$$

$$Sc(R, C) = \frac{n_R}{V(R, C)}$$

∴  $Sc(R, B)$  and  $Sc(R, C)$

$$= \frac{100 \times 100}{100 \times 1} = 1000$$

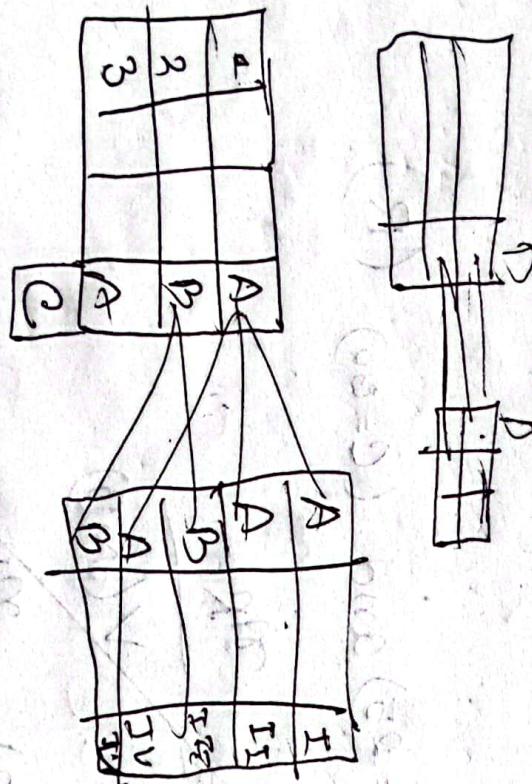
$$Sc = \frac{n_R}{V(R, A)}$$

$$\text{ans} = \frac{Sc}{f_R}$$

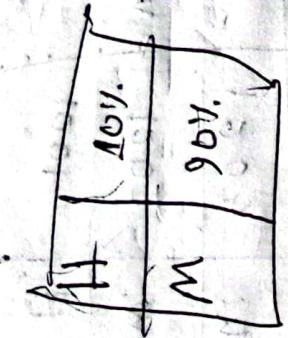
$$b) T(6_{B=25} \& R=30) (R) = \frac{MR}{V(R, B) V(R, r)} = \frac{100}{10 \times 1} = 10$$

$$c) T(6_{B>25}) = \frac{MR}{3} = \frac{100}{3} = 33$$

$$d) T(6_{B>25} \& B=15) (R) = \frac{MR}{V(B>25) \cdot V(B)} = \cancel{\frac{100}{10 \times 3}}$$



A	1201	H	
M	1202	M	
S	1203	M	
F	1204	M	



	Natural			
A	1201	H	H	10%
M	1202	M	H	10%
S	1203	M	H	10%
F	1204	M	H	10%
A	1201	H	M	90%
M	1202	M	M	90%
S	1203	M	M	90%
F	1204	M	M	90%

b)

A	1201	H
M	1202	M
S	1203	M
F	1204	M

Contd

H	100%	Y
M	90%	N
M	90%	Yie

Closed Join

A	1201	H	10%	Y
M	1202	M	90%	N
M	1203	M	90%	Yie
S	1203	M	90%	N
F	1204	M	90%	Yie
F	1204	M	90%	N
			90%	Yie

~~3x2~~      ~~4x3~~      = 6

## Transaction

4/5/22

Abort = undo

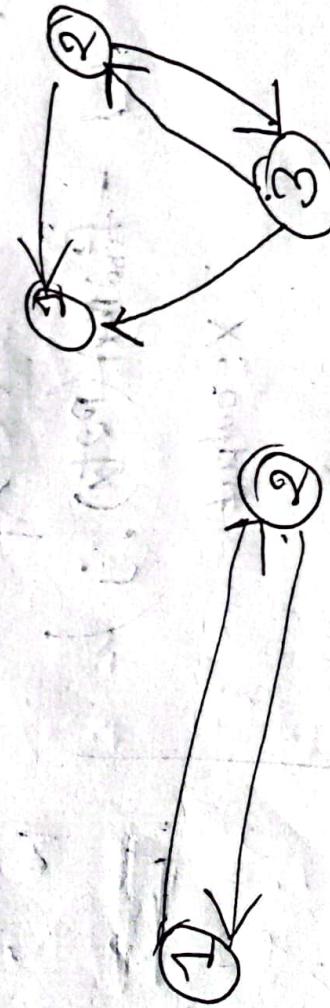
Conflict Serializable: Parallel to serial

or serial: two parallel starts after one other

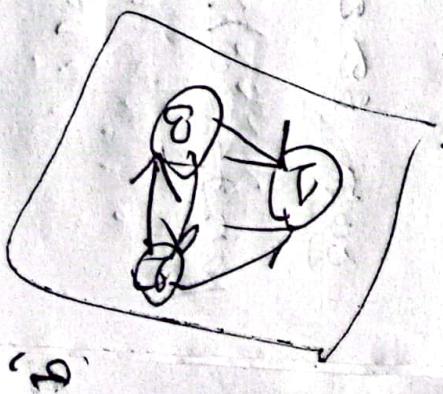
start

• Precedence graph w cycle at vertex

Conflict Serializable.



५०



r2(x)  $\rightarrow$  a b c  
 r1  $\rightarrow$  1 2 3  
 r0  $\rightarrow$  0 0 0

20, 1: 2, 3  
3 : 1, 2

$$\overline{r^2(x) \ r^3(y)} \ \overline{\omega^3(x)}$$

$$p_1(\gamma) \quad w_2(x) \quad w_1(\gamma) \quad \cos(x)$$

Persone (102)

2000-2001  
2001-2002

2

3: for (-) †

$$= \lim_{n \rightarrow \infty} \left( \rho_n(\mu) - \rho_n(\nu) \right)$$

X: output

10

2

200

۲۷

2021-2022, 1 semester

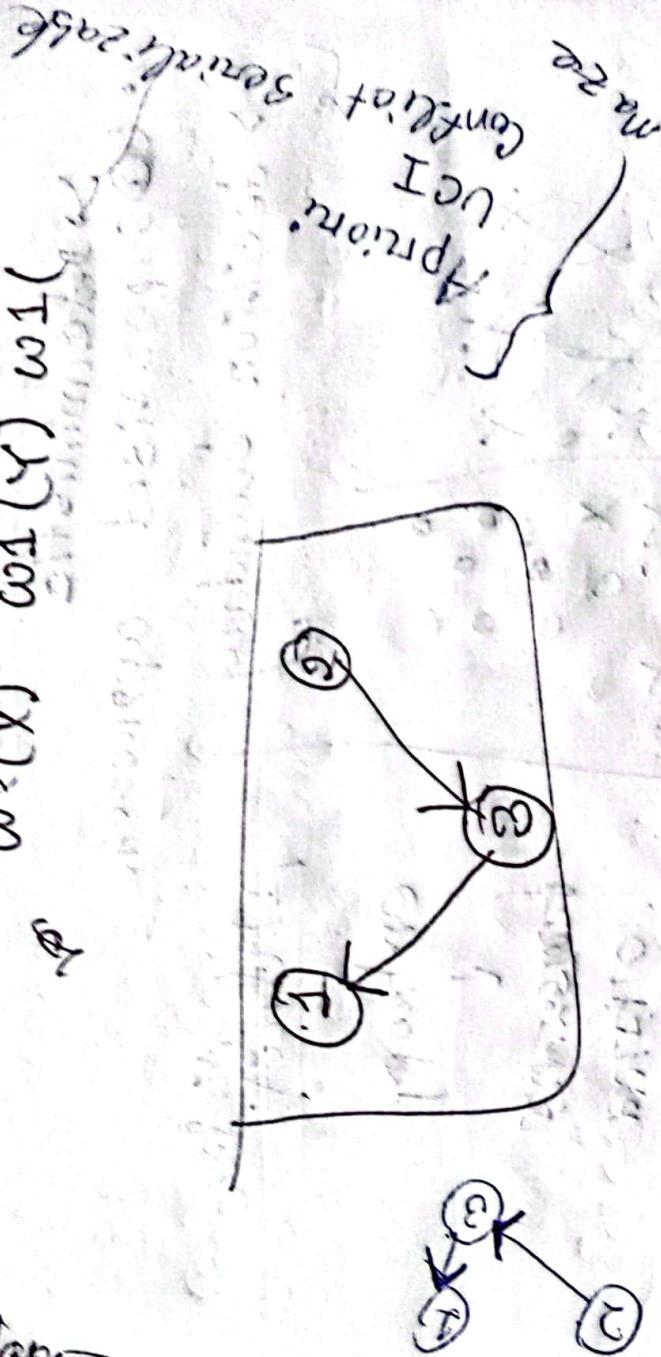
$$P_2(x) P_3(y) w_3(x) w_1(y)$$

$w_1(x)$



$\phi$

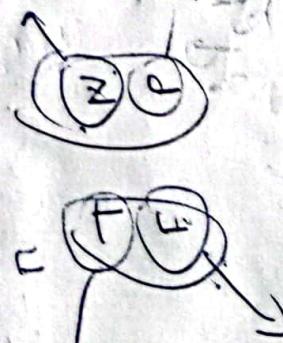
$$w_2(x) w_1(y) w_1(x) w_1(y)$$



non-norm.



Given



$w_2(x)$

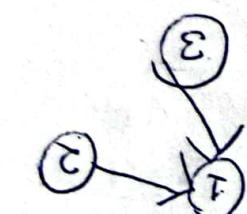
$w_1(y)$

$$= \frac{w_2^T w_1}{d}$$

Non-norm

Posterior

$w_1(y)$

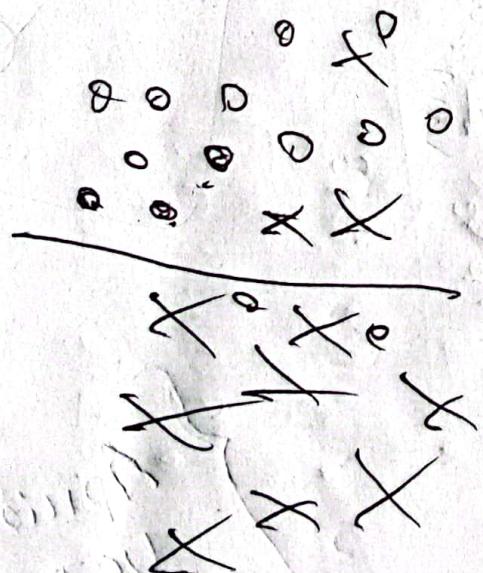


## Decision Tree

- deterministic
- supervised approach.

→ decision boundary

Metric



Entropy

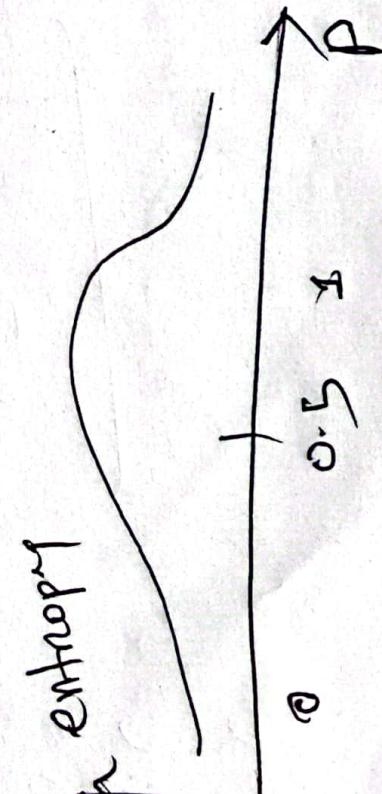
Proportion  
of classes

Entropy

$$\text{Entropy} = \sum_i -P_i \log_2 P_i$$

If  $P_i = 0.5$ , entropy is maximum

Entropy



$\begin{array}{|c|c|} \hline 0 & X \\ \hline 0 & X \\ \hline \end{array}$

50% even, 50% odd  
chaotic scenario. Entropy

is highest.

180° -  
X-axis

-loop

$$E = -\frac{13}{30} \log\left(\frac{13}{30}\right) - \frac{17}{30} \log\left(\frac{17}{30}\right)$$

parent

$$E_{\text{child}} = \text{avg}(E_0, E_1)$$

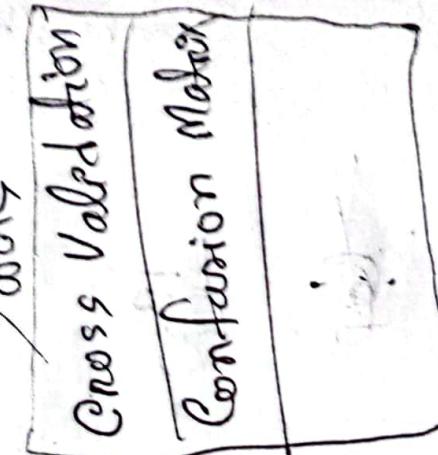
$$E_0 = -\frac{2}{17} \log\left(\frac{2}{17}\right) - \frac{15}{17} \log\left(\frac{15}{17}\right)$$
$$E_1 = -\frac{1}{13} \log\left(\frac{1}{13}\right) - \frac{12}{13} \log\left(\frac{12}{13}\right)$$

$$\text{Info gain} = \text{Gparent} - \text{avg}(Gchild)$$

- # test dataset target score
- BTJ entropy rule

90 / 10

- Sampling



- Pruning
- Index

Task: 5-class assignment (5 classes)

## 510 Parallelism

18/1/22

### Round Robin:

#### Partitioning

##### Point Query

##### Range Query

- Hashed partition
- Round Robin Partition
- Range Partition

### • Skew & binning

### • Parallelism - [Query (Inter query) queries (Intra query)]

⇒ Exam a algorithm (over 01353, \* \* \*)

• Fragment - and - replicate join (over 01351, \* \* \*)

### 3 Partition - Recovery → 4 Parallelism

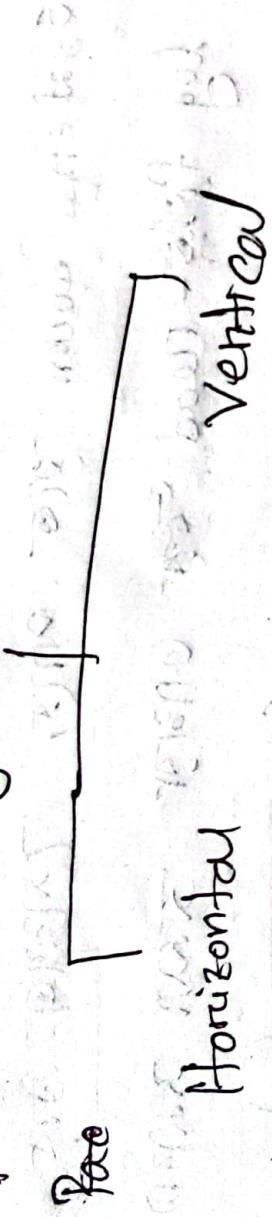
e.g. Join ← e.g. sort (2)

parallel independent

## Distributed Database

24/1/23

- Replication, Segmentation



- Centralized Scheme (name server)
- Naming convention (name + segment + ip)

TC : Transaction  
Coordinator, locally

Tm : Transaction manager, TM

- Two phase commit protocol

Phase 1) Client sends every site  
for comit? (Ready, No, Abort)  
Phase 2) Aborting, (No, Abort)

Phase 2 of the write-ahead log

System writes down the data to disk  
and then writes the log file.

Log file reads the log and then writes it to disk.

Commit T >: Operation Commit message one and  
writer reader reply log - offline one and  
global visibility log

So we just take its result.

Abort T > in Redo:

{Ready T } : Log, System coordinates  
log trace, for local log. We get writer  
for log file.

coordinator, local log fail  
for  $\Delta$

Commit T :  $\Delta$  term

< Abort T > : Undo

< Ready T > : Undo

& if none of the above, then  
wait.

Network Bully algorithm : Coordinator  
will make election after coordinator  
fails

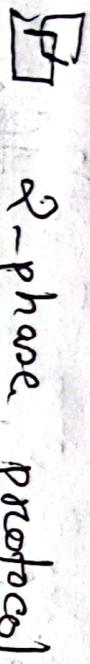
The first node will be selected  
as "Recruiting new  
coordinator" and  
broadcast "Recruiting new  
coordinator" message

and we select the best system (e.g.  
with the most nodes).  
high priority to do our every  
task.

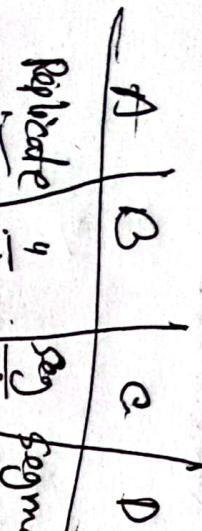
## Join strategy:

$$\text{Dhaka} + \text{Sylhet} = \text{Join}$$
$$5\text{G} + 5\text{G} = 10\text{G}$$

for 5G data transfer it's costly.



• Logfile handed



• Logfile handed

## Function patterns:

4 math + 6 theory  
(24)

Approv: confidence, support count etc.

Naive bayes:

Decision tree: 2 level, 2 level (2A 2B), Entropy

Socialization: control over factors of production

Def: control in transaction

Control	Production	Distribution	Exchange
State	Planning	State	State
Capitalist	Market	Capitalist	Capitalist
Communist	Planning	Communist	Communist
Traditional	Custom	Traditional	Traditional