

SQL Clause

SELECT [DISTINCT] A₁, A₂, ..., A_n
From R₁, R₂, ..., R_n

[] ← Optional clause

[WHERE P]

[GROUP BY ATTRIBUTE [HAVING Condition]]

[ORDER BY Attribute [DESC]]

Query Execution

1. From Clause: It is the first executable clause. It just simply relation or Cross Product of two or more relation.
2. WHERE Clause: It is the second executable clause. It selects the tuple based on specified condition.
3. GROUP BY Clause: It is the third executable clause if used in the query. It groups the table based on the specified attribute.

Q1. Why SQL is most widely used?

Ans. Because Relational Algebra works on 'Set Theory' and SQL works on multi-set (Bag) where duplicate are allowed

e.g.: * set [1, 2, 3, 4]

* multiset [1, 1, 2, 3, 3, 4]

Q2. Commercial Database work on Multiset (Bag) why?

Ans. If we want to perform Union of two table (Relation), SQL just simply one table and simply add the data of another table (∴ it's faster). But in relation we first eliminate duplicates then perform operation (very time consuming).

- * In R.A it eliminated the duplicates so we cannot find the count of multiple tuple value and removing duplicate is expensive
- * In SQL aggregate functions are allowed but not allowed in relational algebra!
- * SQL have lots of features that are not supported by R.A

eg: Supplier (Sid, Sname, Rating)

Parts (Pid, Pname, Color)

Catalog (Sid, Pid, Cost)

WAP to find Sid of Red Color Parts:

→ SELECT C.Sid AS Krishna FROM Catalog AS C, Parts P
WHERE C.Pid = P.pid AND P.color = Red.

Rename

WAP to find Sname of Red Color Parts:

→ SELECT S.Sname FROM Supplier AS S, Part P, Catalog
AS C WHERE P.pid = C.pid AND C.Sid = S.Sid AND P.color
= Red.

Note: Whenever we want multiple copies of the same table to be used then we need to use "Rename/Aliasing".

Aggregate Functions: The functions that take collection of values as input (from set or multiset) and produce single output. It always discards the Null values first.

1. COUNT
2. Sum
3. MIN
4. MAX
5. AVG

eg:

Sid	Branch	Marks
S1	CS	90
S2	JT	70
S3	CS	70
S4	EC	56
S5	CS	Null

Group by
Branch.

Sid	Branch	Marks
S1	CS	90
S3	CS	70
S5	CS	Null
S2	JT	70
S4	EC	56

High frequency

tuple will be

at top

$$\text{COUNT}(\text{marks}) = 4$$

$$\text{Sum}(\text{marks}) = 286 \text{ (Null is removed)}$$

$$\text{COUNT}(\rightarrow) = 5 \text{ (all tuples)}$$

$$\text{Sum}([\text{distinct}] \text{marks}) = 216$$

$$\text{COUNT}([\text{DISTINCT}] \text{marks}) = 3$$

$$\text{AVG}(\text{marks}) = 71.5$$

eg:	A	COUNT(A) : 6
	NULL	COUNT(A) : 0
	NULL	SUM(A) : NULL
	NULL	AVG(A) : NULL
	NULL	MIN(A) : NULL
	NULL	MAX(A) : NULL

Important Points About Aggregate Functions:

- * Arithmetic Operations with NULL give result "NULL".
e.g., $NULL + Any\ Value = NULL$ P.P. Operator

- * Null : Not existed/unknown
 - Two values are not equal
 - Non zero
 - Random ASCII value assigned by DBMS S/W.

- * Aggregate Operator always discard/ignore the Null values.
While Performing the aggregate operation, they first discard or ignore the null values then performed Computations

- * Comparison Operations with NULL give Result "Unknown".

HAVING: Fourth executable clause (if need in query), after "Group By".

- * It is used to select the group which satisfy the conditions (Condition is for each group).

Sl.no	Branch	Mark
S1	CS	60
S2	IT	70
S3	CS	90
S4	IT	60
S5	EC	55
S6	EC	NULL

Sl.no	Branch	Mark
S1	CS	60
S3	CS	90
S2	IT	70
S4	IT	60
S5	EC	55
S6	EC	NULL

Sl.no	Branch	Mark
S1	CS	60
S3	CS	90
S2	IT	70
S4	IT	60

Note: If aggregate operator and other attribute (normal or unaggregate) used in Select Clause is allowed only iff other attribute (non-aggregate attribute) must be in Group By Clause.

SQL SET Operators

1. UNION & UNION ALL
2. INTERSECT & INTERSECT ALL
3. MINUS & MINUS ALL
[EXCEPT] [EXCEPT ALL]

Followed by
Relational Algebra

Not followed
by P.A.

* IP "HAVING" is applied current
"Group By" Clause then Having is
applied to every tuple according
to latest standard of SQL.

eg:

R	S
1	1
1	1
2	2
2	2
4	4
6	5
6	5

$$R \cup S = \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline 4 \\ \hline 6 \\ \hline 5 \\ \hline \end{array}$$

$$R \text{ UNION ALL } S =$$

1
1
2
2
4
6
6
5
5

R INTERSECTS

1
2
4

R INTERSECT ALL S

1
1
2
4

R MINUS S

6

R MINUS ALL S =

1
2
6
6

Other SET Operators:

- ① IN / NOT IN
- ② ANY
- ③ ALL
- ④ EXISTS / NOT EXISTS

Comparison Operators

<, >, <=, >>
not equal

ANY Operator: Compare a value with each value in a set and return true if any value is compared according to given condition.

Nested Queries

① Normal (Independent) Nested Query.

Execution Sequence

Bottom → Top

Inner Query → Outer Query

② Correlated Nested Query.

Execution Sequence:

Top → Bottom → Top

Outer → Inner → Outer
Query Query Query

eg:

Retrieve Sid and marks of the Student who Secured Higher marks?

Query I: Select Sid, marks From Student.

Query II: Select Sid, max(marks) From Student Group By (Sid)

Query III: Select Sid, marks

From Student

✓ Where marks = [Select max(marks)
From Student]

Outer Query

(Inner Query)

Note: Aggregate function Cannot be in lowercase.

ANY Operator: Compare a value with each value in a set and return true if any value is compared according to given Condition

eg: ANY (Op)

$x > \text{ANY}(10, 20, 40)$

O/P ⇒ 11, 12, 13, 14.....

$(x > 10) \text{ or } (x > 20) \text{ or } (x > 40)$

$x < \text{ANY}(10, 20, 40)$

O/P ⇒ 39, 38, 37, 36.....

All Operator: Compare a value with each value in a set and return true if given condition satisfied for each value in the set.

e.g.: $x > \text{All}(10, 20, 30)$ O/p: 31, 32, 33, ...

$x > 10 \text{ AND } x > 20 \text{ AND } x > 30$

Q1. Find the name of Supplier whose turnover is better than the turnover of any (Some) Supplier of Delhi.

SUPPLIER

Sno.	Sname	City	turnover
1	A	Delhi	4 Cr.
2	B	Bengaluru	5 Cr.
3	C	Delhi	6 Cr.
4	D	Kochi	7 Cr.

SELECT Sname From Supplier

WHERE City >= Delhi AND

turnover > ANY (Select turnover From Supplier WHERE City = Delhi)

O/p: Any B ALL D

Q2. Retrieve eid who get more salary than any employee of dept = 5?

Emp	Eid	Dno	Salary
E1	5	50K	
E2	3	20K	
E3	5	30K	
E4	3	40K	
E5	4.	60K.	

Query I: SELECT Eid From Emp WHERE

$dno >= '5'$ AND Salary > ANY

(SELECT Salary From Emp WHERE dno = '5')

Query II : (Without ANY Operator)

SELECT Eid From Emp WHERE

$dno >= 5$ AND Salary > (Select min(Salary) WHERE dno = 5).

O/p: Any E4 ALL E5

E5

Q3. Consider the following relation Cinema (theatre, address, capacity)
Which of the following option can be at the end of the SQL Query
"SELECT P1.address From Cinema P1" such that it always finds the
address of theatre with maximum capacity?

WHERE P1.capacity >= All (Select P2.capacity From Cinema P2).

Ans

Supplier (Sid, Sname, Rating)

Paint (Pid, Pname, Color)

Catalog (Sid, Pid, Cost)

Supplier

Sid	Sname	Rating
S1	Abhay	8
S2	Praveen	9
S3	Priya	9
S4	Pooja	9
S5	Alok	9

Catalog

Sid	Pid	Cost
S1	P1	8K
S2	P2	7K
S3	P3	8K
S4	P4	9K

Paint

Pid	Pname	Color
P1	A	Red
P2	B	Green
P3	C	Blue
P4	D	Yellow

1. W.A.Q to find Sid of Supplier whose rating > 8

→ SELECT Sid FROM Supplier WHERE Rating > 8 OR

SELECT T.Sid FROM Supplier T WHERE T.Rating > 8

S2
S3
S4
S5

2. W.A.Q to find Sid whose parts cost greater than 7K.

→ SELECT X.Sid FROM Catalog X WHERE X.Cost > 7000

OR SELECT Sid FROM Catalog WHERE Cost > 7000

S2
S4

3. W.A.Q to find Pid of Red color paint

→ SELECT Pid FROM Paint WHERE Color = 'Red'

P1
P3

4. W.A.Q to find Sid of Supplier who Supplied Red Color Paint.

→ SELECT C.Sid FROM Catalog C, Paint P WHERE C.Pid = P.pid
AND P.color = 'Red'.

S1
S2

5. W.A.Q to find 'Sname' who Supplied Red Color Paint.

→ SELECT S.Sname

From Catalog C, Paint P, Supplier As S

WHERE C.Pid = P.pid AND

C.Sid = S.Sid AND

P.color = 'Red'.

Sname
Abhay
Praveen

Q. W.A.Q to Find Sid who not supplied Red Color Part?

→ SELECT C.Sid From Catalog C, Part P WHERE C.pid = P.pid AND P.color <> 'Red'.

S1
S4

Note:

$$x = 5$$

$$R[1, 2, 3, 4, 5]$$

$$x \text{ IN } R$$

$$'5' \text{ IN } [1, 2, 3, 4, 5] \Rightarrow \text{True}$$

$$x \text{ NOT IN } R$$

$$'5' \text{ NOT IN } [1, 2, 3, 4, 5]$$

$$\Rightarrow \text{False}$$

Query: SELECT Sid From Catalog WHERE pid IN (SELECT pid From Parts WHERE color = Red).

Using "IN"

(i) P₁ IN (P₁, P₃): True

↳ Then S₁ Selected

(ii) P₂ IN (P₁, P₃): False

↳ Select nothing

(iii) P₃ IN (P₁, P₃): True

↳ P₃ S₂ Selected

(iv) P₄ IN (P₁, P₃): False

↳ Nothing Selected.

Sid
S1
S2

Using "NOT IN"

P₁ NOT IN (P₁, P₃): False

P₂ NOT IN (P₁, P₃): True

↳ S₁ Selected

P₃ NOT IN (P₁, P₃): False

P₄ NOT IN (P₁, P₃): True

↳ S₄ Selected

Sid
S1
S4

← Color nor Red

Homework: (Using the Supplier, Parts, Catalog table)

- W.A.Q to find Sid who Supplied Some Red Color parts or Some Green Color parts.
- W.A.Q to find Sid who Supplied Some Red Color part and Some Green Color Part.

STUDENT (Id, Name, Branch, Marks)

- W.A.Q to Retrive Id of Student who Secured Higher marks (First highest marks) in → R.A, SQL, max (Aggregate) Function.
- W.A.Q to Retrive Id of Student who Secured 2nd highest marks. → (i) R.A (ii) SQL (iii) Max (Aggregate) Function.

Supplier (Sid, Sname, Rating) }
 Parts (Pid, Pname, Color) } Relations
 Catalog (Sid, Pid, Cost)

- (i) W.A.Q to find Sid's who supplied some red color parts or some green color parts

RA: $\Pi_{\text{Sid}} \left[\begin{array}{l} \sigma_{\text{C.pid} = \text{P.pid}} \wedge (\text{Catalog} \times \text{Parts}) \\ \text{P.color} = \text{Red} \vee \text{Green} \\ \text{P.color} = \text{Red} \text{ OR } \text{P.color} = \text{Green} \end{array} \right]$

SQL: `SELECT Sid`

`FROM Catalog C, Parts P WHERE C.pid = P.pid AND`
`P.color = 'Red' OR 'Green'.`

- (ii) W.A.Q to find Sid's who supplied some 'Red' color part and some 'Green' color part.

RA: $\Pi_{\text{Sid}} \left[\begin{array}{l} \sigma_{\text{C.pid} = \text{P.pid}} \wedge (\text{Catalog} \times \text{Parts}) \\ \text{P.color} = \text{Red} \end{array} \right] \cap \Pi_{\text{Sid}} \left[\begin{array}{l} \sigma_{\text{P.pid} = \text{C.pid}} \wedge \\ \text{P.color} = \text{Green} \end{array} \right]$

SQL: `Select C.Sid From Catalog C, Parts P WHERE`
`C.pid = P.pid AND P.color = 'Red'`

INTERSECT

`Select C.Sid From Catalog C, Parts P WHERE`
`C.pid = P.pid AND P.color = 'Green'.`

STUDENT

Sid	Name	Branch	Mark
S1	Khubbi	CS	80
S2	Qwarab	ME	72
S3	Renchir	CS	92
S4	Bhanya	IT	90

(i) W.A.Q to find unique Sid of student who secured higher (1st highest mark).mark.

O/P:

Sid
S3

RA: $\Pi_{\text{Sid}}(\text{Student}) - \Pi_{\text{Sid}} \left[\text{Student} \times \min_{\text{mark} < m} \right]$

$\sigma_{T, N, B, M}(\text{Student})$

SQL: $\text{SELECT Sid From STUDENT}$
EXCEPT

$\text{SELECT T1.Sid From STUDENT as T1, STUDENT as T2 WHERE}$
 $T1.marks < T2.marks$

Using MAX Function: $\text{SELECT Sid From STUDENT WHERE}$

$\text{marks} = \left[\text{SELECT max(marks) From STUDENT} \right]$

(ii) N.A.Q to retrieve Sid of Student who Secured 2nd highest marks: $Op = S4$

Rn: $P(Temp, T)$ $\left[\text{STUDENT} \bowtie_{\text{Sid marks}} P_{IN, B, M}(\text{STUDENT}) \right]$

$T_{Sid}(Temp) = T_{Sid} \left[\text{Temp} \bowtie_{\text{marks} < m} P_{IN}(Temp) \right]$

SQL:

WITH Temp(Sid, marks)

Used to sub query or as a database
 $\text{SELECT T1.Sid T1.marks From STUDENT T1, STUDENT T2}$
 $\text{WHERE T1.marks} < \text{T2.marks}$

$\text{SELECT Sid From Temp}$

EXCEPT

$\text{SELECT T1.sid From Temp T1, Temp T2}$
 $\text{WHERE T1.marks} < \text{T2.marks}$

Using MAX Function:

SELECT Sid

$\text{From STUDENT WHERE Tmarks} =$

$\left[\text{SELECT max(marks)}$

From STUDENT

$\text{WHERE marks} <$

$\left[\text{Select max(marks)}$
 $\text{From STUDENT} \right]$

Q1. The following relation records the age of 500 employee of a company where empNo { indicating the employee number } is the key.

empAge (empNo, Age)

Consider the following R.n expression:

T empNo (empAge Δ (age > age)) P empNo, age (empAge)

What does the above expression generate?

Anc Employee numbers of all employees whose age is not the minimum.

Q2. Retrieve Sid and Marks of the Student who secured highest marks

\times Query I: Select Sid, max(marks) { Syntax error }

From Student

\times Query II: Select Sid, max(marks) { Syntax correct but it gives all Sid and their marks, not giving topper id and marks }

From Student

Group By (Sid)

✓ Query III: Select Sid, marks
From Student

Where marks = (Select max(Marks) From Student)

Note: * When aggregate operator and other operators attribute (non aggregate) used in the select clause is allowed only if other attribute must be in the group by clause.

* Aggregate function cannot be in lower clause.

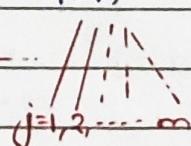
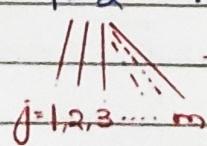
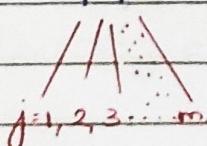
Correlated Nested Query Execution : Inner Query Using attribute defined in outer query

for(j=1; j<=m; +j)

i=1

i=2

i=n



eg: of Correlated nested query:

SELECT C.sid

From Catalog C

WHERE EXISTS

SELECT * From Parts P

WHERE P.pid = C.pid

AND Color = 'Red'

Toner Query Using correlated

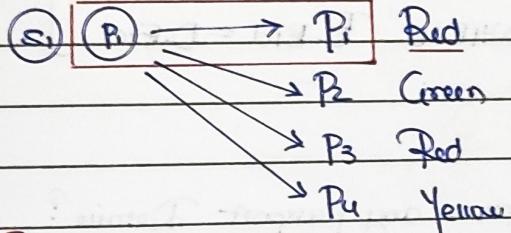
Defined in Outer query

EXISTS (Check): Return True if Toner Query Result Not Empty

NOT EXISTS: Return True if Toner Result Empty.

Execution of above query:

1st Iteration



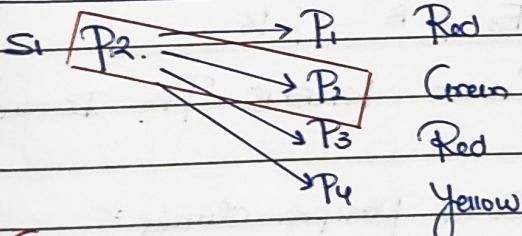
Catalog

Sid	Pid
S ₁	P ₁
S ₂	P ₂
S ₃	P ₃
S ₄	P ₄

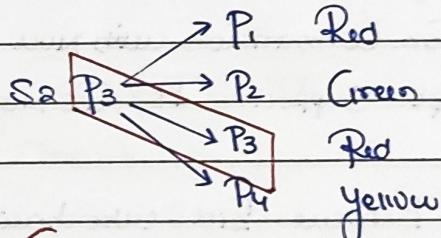
Parts

Pid	Color
P ₁	Red
P ₂	Green
P ₃	Red
P ₄	Yellow

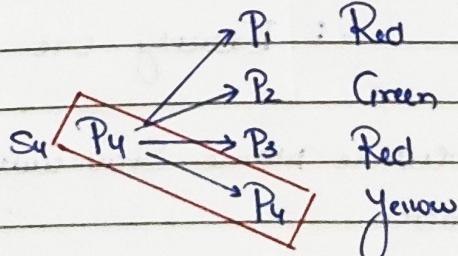
2nd Iteration



3rd Iteration



4th Iteration



Exist:	S_1	NOT Exist:	S_1
Color = 'Red'	S_2	Color != 'Red'	S_4

- Before Exist and NOT Exist no attribute is required.
- Before IN & NOT IN Attribute is required.

Q3 Given relational Schema Emp (Eid, Ename, Salary)
Department (Eid, Dname, Code)

Retrieve Employee Id who have no Department.

✗ Query I: Select Eid from Emp E, Dep D
where E.Eid <> D.eid

✓ Query II: SELECT Eid From Emp E WHERE
NOT EXISTS [SELECT * From Dep D
WHERE E.Eid = D.Eid]

Ans Q2 is correct & Q1 is wrong

Q4 Retrieve Eid who does not have any passport Details?

→ SELECT Eid

From Emp

WHERE Pno IS NULL

eid	ename	Pno
E1	A	20
E2	B	NULL
E3	B	16

* For Comparision with null

SQL Support Is/Is NOT Clause WHERE Pno IS NULL

(ii) Retrieve Eid who having Some Passport Details:

→ SELECT Eid From Emp WHERE Pno IS NOT NULL

Regular Expression: % : Zero or more Character
— : Exactly one Character.

Q5 Retrieve Student whose Name Starting 'S' & end with M
& at least 5 Characters.

Ans S — % M.

→ SELECT * From Student WHERE Sname LIKE '% M'

LIKE: It is used to compare to specify certain search conditions for a pattern in a column.

Q6. Retrieve all Students whose name not starts with 'C'?

Ans SELECT * From Student WHERE Sname NOT LIKE 'C%'.

eg: SELECT * From Student WHERE Sname

LIKE A% → Starts with A

%.J → ends with J

%.I%. → Contains I

'___' → All 4 length name

'S__' → Starts with S and exactly 4 length

'S___%' → Starts with S and atleast length is 4

1. R x S : Cross Join

2. R US : Union Join

3. R RS : Inner Join

4. R RT S : Outer Join

ORDER BY : Orderby clause is used to

Sort the result

* By default : Ascending [ASC]

Descending [DESC]

Q7. Consider the relations : Supply (SupplierId, ItemCode) with 100 tuples
 Inventory (ItemCode, Color) with 200 tuples
 Let P and Q be the number of maximum and minimum records in Supply Join inventory, the value of P+Q is _____. (Itemcode is Primary key of Supply table).

Ans Supply : Preferred Relation
 F.K : Item code

maximum no. of Tuples : 100

minimum no. of Tuples : 100

Inventory : Preferred Relation

Here Itemcode is key so

Itemcode cannot be NULL

So min tuple = 100

Join using foreign key concept.

(ii) $R(ABC)$ $S(BCD)$

B is foreign key

→ 100 tuple

S (BCD)

↓
500 tuple

Maximum Tuples in $R \times S = 1000$

(iii) $R(ABC)$ $S(BCD)$

B is foreign key

→ 100 tuple

S (BCD)

↓
400 tuple

Maximum Tuples in $R \times S = 1000$

Note: * The value present in foreign key must be present in primary key of referenced relation.

* Foreign key may contain "duplicate" and "null" values

* Table which contains \rightarrow Referencing Relations (Table)
foreign key

* Whenever Two Table (Relations) are Joined (natural join) with respect to Primary key and foreign key then maximum number of tuples in the resulting relations is equal to number of tuple in the referencing relations.

Q8. Consider the relations:

Employee (Eid , $Ename$) with 200-tuple

Department (Eid , $Dname$, $Color$) with 100-tuple.

Let p and q be the number of maximum and minimum records in Employee Join Department. The value of $p+q$ is 200. (Eid is FK in Department table).

Ans

Maximum = 100 } 200 Total

Minimum = 100 }

Q9.

Consider the relations:

$R_1(A, B, C)$ with n tuples.

$R_2(B, D, E)$ with m tuples (B is FK in R_1)

Maximum tuples in $R_1 \times R_2 = n$

Minimum tuples in $R_1 \times R_2 = 0$ (If B is foreign key)

Minimum tuples = 0 (B is not FK,

B is not key, or B is having Null value)

Q10.

Consider the relations:

$R_1(A, B, C)$ with n tuples

$R_2(B, D, E)$ with m tuples (B is not FK in R_1)

Maximum tuples in $R_1 \times R_2 = n$

Minimum tuples in $R_1 \times R_2 = 0$

Q11. Consider the join $R \times S$ between $R(ABC)$ and $S(ADE)$ with attribute

A being the primary key in both relations and attribute A in S is a Foreign key referring to attribute A in R . R has n tuples and S has m tuples, maximum and minimum tuples in $R \times S$

respectively are?

$R(ABC)$ min (m, n), max (m, n)

Referenced Relation			Referring Relation		
A	B	C	A	D	E
1			1		
2			2		
3			4		
4			5		
5			6		
6			7		
7			8		

$A : (\text{unique} + \text{not null})$

Note:

$R(ABC)$

Referencing Relations.

$S(ADE)$

'FK & also primary key.'

Here, the value of A in Relation S contain not only {Null}, but not duplicates also.

A in S : Must be unique + not null

But those value of ' A '

must be present in ' A ' of Relation ' R '. (Foreign Key Concept)

Q12. Consider the relations:

$R_1(A, B, C)$ with 3 tuples

$R_2(D, B, E)$ with 3 tuples $\{A \rightarrow B, B \rightarrow C, D \rightarrow B, B \rightarrow E\}$

Maximum and Minimum tuple in $R_1 \bowtie R_2$.

Ans

$R_1(ABC)$

$R_2(DBE)$

$[A]^+ = [ABC]$

$[D]^+ = [DBE]$ D: Primary Key

A is P.K.

$\rightarrow B$ is non key attribute in both relation R_1 & R_2

Maximum = ABC

Minimum = O

Q13. The following functional dependencies hold for relations $R(A, B, C)$ and $S(B, D, E)$ Fd for both the tables. $\{B \rightarrow A, A \rightarrow C\}$

The relation R contains 200 tuples and the relation S contains 100 tuples. What is the maximum number of tuples possible in natural join $R \bowtie S$?

Ans

$R(A \leftarrow BC)$

$S(BDE)$

B: Primary

Referencing

Min no of tuples = 100

(Referred Relations) $\xrightarrow{\text{Key}}$ Relation

Q14. Consider two relations $R_1(A, B)$ with the tuples (1,5), (3,7) and $R_2(A, C) = (1,7) (4,9)$

Assume that $R(A, B, C)$ is the full natural outer join of R_1 and R_2 .

Consider the following tuples of the form (A, B, C) ; $a = (1, 5, null)$,

$b = (1, null, 7)$, $c = (3, null, 9)$, $d = (4, 7, null)$, $e = (1, 5, 7)$

$f = (3, 7, null)$, $g = (4, 9, null)$. Which of the following statements is correct?

Ans.

Relations contain e,f,g

but not a,b

$R \bowtie S$

$R \bowtie S$

A	B	C
1	5	7
3	7	N
4	N	9

e

f

g

Q15

Suppose $R_1(A, B)$ and $R_2(C, D)$ are two relation schemas. Let r_1 and r_2 be the corresponding relation instances. B is a foreign key that refers to C in R_2 . If data in r_1 and r_2 satisfy referential integrity constraints, which of the following will always be true?

Ans

$$\text{Fl}_B(r_1) = \text{Fl}_C(r_2) = \emptyset$$

$$[1, 2] - [1, 2, 3, 4] = \emptyset$$

R_1	\xrightarrow{Fk}	R_2	
A	B	C	D
1		1	
1		2	
2		3	
2		4	

Q16.

Consider the following Database Schema named water_schemas : the no. of tuples returned by the following SQL Query is :

WITH total(name, capacity) as
SELECT dist_name, sum(capacity)
from water_schemas

Group By dist_name

WITH total_avg(capacity) as

SELECT avg(capacity)

From total

SELECT name

From total, total_avg

WHERE total_capacity >= total_avg_capacity.

Scheme-no	dist_name	Capacity
1	Ajmer	20
1	Bikaner	10
2	Bikaner	10
3	Bikaner	20
1	Churu	10
2	Churu	20
1	Dungarpur	10

Ans

Total

Name	Capacity
Ajmer	20
Bikaner	40
Churu	30
Dungarpur	10

Total-Avg

Avg Capacity = 25

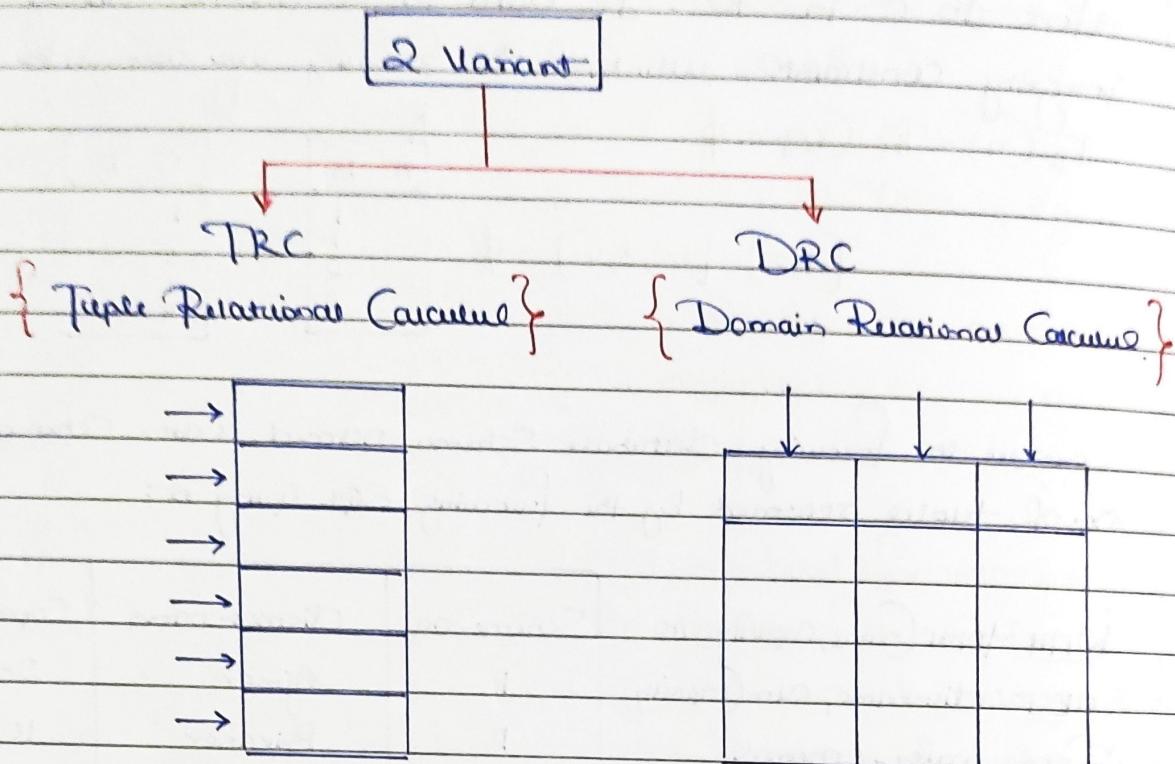
Total Cap > Total Avg Capacity

Name
Bikaner
Churu

Ans = 2 Tuple

Relational Calculus (Based on Predicate Calculus)

Non Procedural query language and has



Tuple as a basic unit.

DRC Takes a column
(Domain) and examine it.

TRC

Query describe result in the form of set of tuples.

Query describe result in the form of set of columns (Domain).

Tuple Variable

It is a variable that takes on tuples of a relation schema as value.

Domain Variable

It is a variable that range over the domains of some attribute (column).

Form of Query: { T | P(T) }.

↑
Tuple
variable

Formula which describes
Tuple Variable

Form of Query: $\{ \langle x_1, x_2, x_3, \dots, x_n \rangle \mid P(x_1, x_2, x_3, \dots, x_n) \}$

↑
Domain
Variable

↑
Formula which describe
Domain Variable

TRC

$\{ T \mid P(T) \}$

T : Tuple variable

$P(T)$: Formula over Tuple(T)
Such that $P(T)$ is satisfied

Form of Query

$T \{ \text{SELECT } \dots \}$

$P(T) \{ \text{From } \dots \text{ WHERE } * \text{ Condition } *$

Note: Formula of TRC is expressed using first order logic.
(predicate logic)

First Order Logic:

Variable are of 2 Type

① Free Variable

② Bounded Variable

Quantifiers:

1. \forall : For all (All)

2. \exists : There Exists (Any)

Bounded Variable: If tuple variable is preceded by the quantifier then it is bounded variable.

e.g.: Belong

1. $\forall E$ Supplier

2. $\exists E$ Supplier

Free Variable: If tuple variable is not bounded by the quantifier then its free variable.

Note: The result of TRC should be free variable.

Q1. Loan (Loanno, Branch name, Amount)

Find Loanno of amount above 50000.

Ane.

TRC:

$$\{ T \mid T \in \text{Loan} (\text{Total amount} > 50,000) \}$$

DRC:

$$\{ \langle \text{Loanno, Branchname, Amount} \rangle \mid \langle \text{Loanno, Branch, amount} \rangle \in \text{Loan} (\text{Amount} > 50000) \}$$

Q2. Retrieve Sid of the Supplier who supplied some Red color parts.

Ane.

Relational Algebra:

$$\pi_{\text{Sid}} (\sigma_{\text{C.pid} = \text{P.pid}} (\text{Catalog} \times \sigma_{\text{Color} = \text{Red}} (\text{Parts})))$$

TRC:

$$(T \mid \exists C \in \text{Catalog} \exists P \in \text{parts} (P.\text{color} = \text{Red} \wedge P.\text{pid} = C.\text{pid}) \wedge T = C.\text{sid}))$$

TRC:

→	
→	
→	
→	
→	
→	

→ Take a Tuple and examine it

→ Tuple as a basic unit.

DRC : → Take a Column

→ Range over the column or Domain & examine it.

TRC : Unsafe Operator Occur

Unsafe Operations : Student (t)

$\neg(t)$: Not belongs to student table (Infinite/ Universe all tuple).

So it is Unsafe operation.

FILE ORGANIZATION & INDEXING



- * Database is a collection of files
- * Each file is a collection of records
- * Each record is a sequence of fields
- * DB is divided into number of blocks
- * Each record block is divided into records
- * Record can be stored in a block

Blocking Factor [BF] : Average number of records per block.

$$\text{Blocking Factor} = \frac{\text{Block Size}}{\text{Record Size}}$$

{ Block → Record
Size Size }
→ always greater.

Organization

① Spanned Organization.

→ A record can be stored more than one block.

(Partially can be stored)

e.g. Block size = 100B

Record size = 40B

$$\text{Blocking Factor} = \frac{100B}{40B}$$

= 2.5 records/block.

② Unspanned Strategy Organization.

→ Record can be stored/belong to a particular block.

(Cannot stored partially)

e.g. Block size = 100B

Record size = 40B

$$\text{Blocking Factor} = \left\lceil \frac{100B}{40B} \right\rceil = [2.5]$$

= 2 records/block

	R ₁ 40B
Block 1	R ₂ 40B
	1/2 R ₃ 20B
	1/2 R ₃ 20B
Block 2	R ₄ 40B
	R ₅ 40B

	R ₁ 40B
Block 1	R ₂ 40B
	20B
	R ₃ 40B
Block 2	R ₄ 40B
	20B

Spanned Organisation:

- * Advantage: no memory wastage
- * Disadvantage: block access
Cost: Increased
- * Suitable for variable length record.

Unspanned Organisation:

- * Advantage: Block access reduced, easy manage
- * Disadvantage: Wastage of memory
(Internal fragmentation)
- * Suitable for fixed length record.

Note: In Spanned organisations no memory is wasted but I/O cost is more (block access increase)

- * But in Unspanned Organisations memory is wasted but input-output cost is less compared to Spanned organisations.
- * Default organisation is Unspanned organisations.
- * I/O Cost: Input output cost means number of blocks transferred from a secondary memory to main memory in order to access some records.
- * Search Key: Attribute used to access the data from DB organisation of records is :
 - ① Ordered file organisation.
 - ② Unordered file organisation.

File Organisations:

Data Records can be : ① Fixed length records

② Variable length records

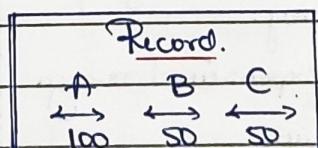
① Fixed Length records:

Create table R
(A char (100),

B char (50),

C char (50)) ;

Block.



R ₁	200B
R ₂	200B
R ₃	200B
⋮	

Block leaders (Used to store offset table used to access records).

② Variable Length Records

Create table S

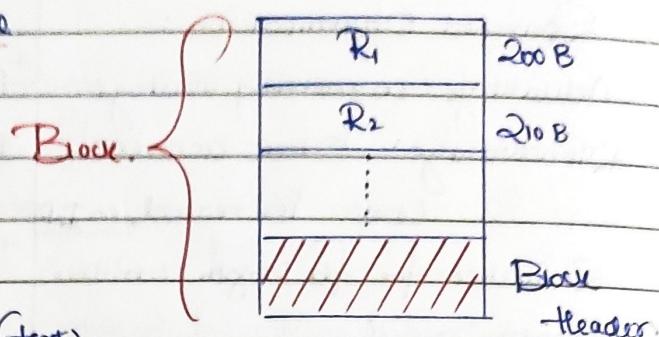
(D char(100),

E char (50),

F text);

$$\therefore \text{Size}(D) + \text{size}(E) + \text{size}(F \text{text})$$

$$100 + 50 \quad * \text{not fixed}$$



DB file with all records fixed length

DB file with variable length record

Both are possible,

DBms

Block Factor: Maximum possible records per block.

Block size : B bytes

Block header size = H bytes

Record size : R bytes

→ Block factor for unspanned : $\left\lfloor \frac{B-H}{R} \right\rfloor$ record/block. because for

fixed length record unspanned is preferred.

→ Block factor for spanned : $\frac{B-H}{R}$ record/block.

	Ordered File	Unordered File
1.	Searching easy.	Insertion easy
2.	Insertion expensive (re-organizing the file)	Searching expensive. <u>Linear Search</u>
3.	<u>Binary Search</u> To access a record avg no. of blocks = $\lceil \log_2 B \rceil$ B: block data	Average case = $\frac{B}{2}$ Worst Case = B.

File of Ordered records (Sorted File)

Ordered (Sequential) file :

- * Records sorted by ordering field
- * Can use ordering key if ordering field is a key field.

Advantages:

- * Reading records in order of ordering key value is extremely efficient.
- * Finding next record
- * Binary Search technique (to access a record avg no of block access = $\log_2 B$)

File of Unordered Records (Heap File)

Heap (or pile) file :

- * Records placed in file in order of insertion.

- * Inserting a new record is very efficient

- * Searching for a record required Linear Search (avg. no of block access).

- * Deletion techniques : 1. Rewrite the Block 2. Use deletion marker.

eg: Suppose that record size $R = 150$ bytes, block size $B = 512$ bytes,

$T = 30000$ records.

then we get blocking factor :
$$\left[\frac{\text{Blocking size}}{\text{Record Size}} \right] = \left[\frac{512B}{150B} \right]$$

$$= 3 \text{ Records/Block}$$

No of file blocks $b = \frac{30,000}{3} = 10,000$ data block.

Indexing:

- * Index are used to improve the search efficiency.

- * To reduce I/O cost.

- * Provide a faster way to access the data.

- * Index is an ordered file.

- * Search key: Attribut to set of attributes used to look up records in a file.



→ Points to a place where key is available.

- * An index file consists of records (called index entries) of the form shown above.

$$\text{One Index Record} = \frac{\text{Size}}{\text{Size}} + \frac{\text{Key}}{\text{Key}} + \frac{\text{Pointer}}{\text{Pointer}}$$

- * Index files are typically much smaller than the original files.

- * Two basic kinds of indices:

1. Ordered indices: Search keys are stored in sorted order
2. Hash indices: Search keys are distributed uniformly across buckets using a hash function.

- * Index file block size is same as DB file block size.

- * Block size of index file = Block size of DB file

- * One index record size = Search key size + size of block pointer.

Note: To access a record average number of block access
 $= \log_2 B + 1$.

↓ →
 Index block Data block
 access access [B: Index Block]

Q1. Given the following data file

EMPLOYEE (Name, Ssn, Address, Job, Sal...)

Suppose that:

- * record size R = 150 bytes, block size B = 512 bytes, $\gamma = 300$ records
- * For an index on the Ssn field, assume that the field size USN = 9 bytes, assume the record pointer size PR = 7 bytes
- * Then

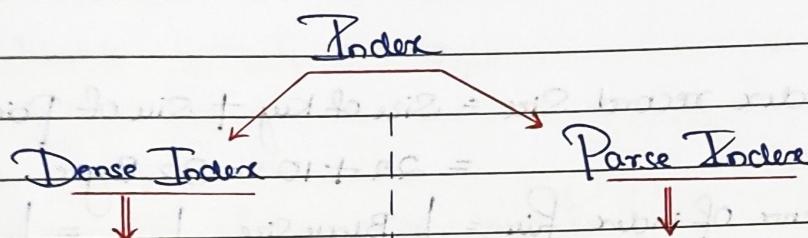
Ans.

One index record Size = $9+7 = 16$ Byte

Block factor of index file = $\left[\frac{512B}{16B} \right] = 32$ Index record Per Block

Total no. of index entries = 30,000, then total no. of index blocks = $\left[\frac{30,000}{32} \right] = 938$ Index blocks

Average no. of block access = $\log_2 B_i + 1$
 (To access a record). = $\lceil \log_2 938 \rceil + 1 = 10 + 1 = 11$ Block access



- * Index entry created for every search key value.
- * Number of index entries = Number of DB records.
- * Index entry is created for some search key values.
- * Number of index entries = Number of blocks.
- * Applicable when records are sequentially ordered on search key.

In Sparse Index file:

To locate a record with Search-key value k we:

- * Find index record with largest search-key value $\leq k$
- * Search file sequentially starting at the record to which the index record points.

Q2 No. of records of file = 16384 (2^14). Block size = 4096 Bytes.
 Record size = 256 Bytes. Search key size = 22 Bytes.

Pointers : 10 Bytes

- (i) P10 Cost to access record without index based on
- Ordered file = _____
 - Unordered file = _____

- (ii) Index dense [using] : (a) # of index blocks at 1st level =
 (b) I/O cost to access record : [using 1st level] =

(iii) By using sparse index :

(a) # of index blocks at 1st level : _____

(b) I/O cost to access record : [using 1st level] = _____

Ans:

$$\text{Block factor of DB file} = \left\lceil \frac{\text{Block size}}{\text{Record size}} \right\rceil = \left\lceil \frac{1096 \text{ B}}{256 \text{ B}} \right\rceil = 2^7 \text{ Record per block}$$

$$\text{Total # data blocks} = \left\lceil \frac{\# \text{records}}{\text{BF of DB file}} \right\rceil = \left\lceil \frac{16,384}{16} \right\rceil = 2^{10} = 1024 \text{ data blocks}$$

Indexing:

$$\text{One index record size} = \text{Size of key} + \text{Size of pointer} \\ = 22 + 10 = 32 \text{ Bytes}$$

$$\text{Block factor of index file} = \left\lceil \frac{\text{Block size}}{\text{Index record size}} \right\rceil = \left\lceil \frac{1096 \text{ B}}{32 \text{ B}} \right\rceil = 2^7 \text{ Index entries / block}$$

Without index:

Ordered file:

$$\text{To access a recg a no. of block access} = \lceil \log_2 B \rceil \\ = \lceil \log_2 1024 \rceil \\ = 10 \text{ Block Access}$$

Unordered file:

$$\text{To access a record average of block access} = \frac{B}{2} = \frac{1024}{2} = 512 \text{ Block access}$$

With Index: Dense.

$$\text{Total no. of index entries} = 16,384 [2^4] \\ (\# DB records)$$

$$B_{fi} \text{ (Block factor of index file)} = \frac{128}{\text{Index entries per block}}$$

$$\text{Total no. of index blocks} = \frac{16,384}{128} = 2^7 \text{ Index blocks}$$

Sparse

$$\text{Total # of index entries} = 1024 \\ (\# DB records)$$

$$\text{Block factor of index file} = \frac{128}{1024} \text{ Index entries per block}$$

$$\text{Total no. of index blocks} = \frac{1024}{128} = 2^3 \text{ Index blocks}$$

With Index: Dense

To access a record index avg # block access = $\lceil \log_2 B \rceil$
 $= \lceil \log_2 128 \rceil = 7$.

To access a record using index average # of block access = $\lceil \log_2 128 \rceil + 1$
 $= 8$ Ans

Sparse

To access index avg # block access
 $= \lceil \log_2 8 \rceil = 3$.

To access a record using index avg number of block access = $\lceil \log_2 B_i \rceil + 1$
 $= \log_2 8 + 1 = 4$ Ans

Meta Data (Data Dictionary)

- record format
- field format
- number of bytes allocated for file
- number of records in file.
- ordered field of file, etc.

Index File:

$$BF \text{ (Block Factor) of Index} = \left\lfloor \frac{B-1}{K+p} \right\rfloor \text{ entries/block}$$

Index file	K_1	P_1
	K_2	P_2
	⋮	⋮

$$\therefore BF = \left\lfloor \frac{B-1}{K+p} \right\rfloor$$

{# of DB file blocks}

$$\left[\frac{B-1}{R} \right] \text{ record/block} \times \left[\frac{B-1}{K+p} \right] \text{ entries/block}$$

{# of index blocks}

Type of IndexSingle Level Ordered Index

1. Primary index
2. Clustering index
3. Secondary index

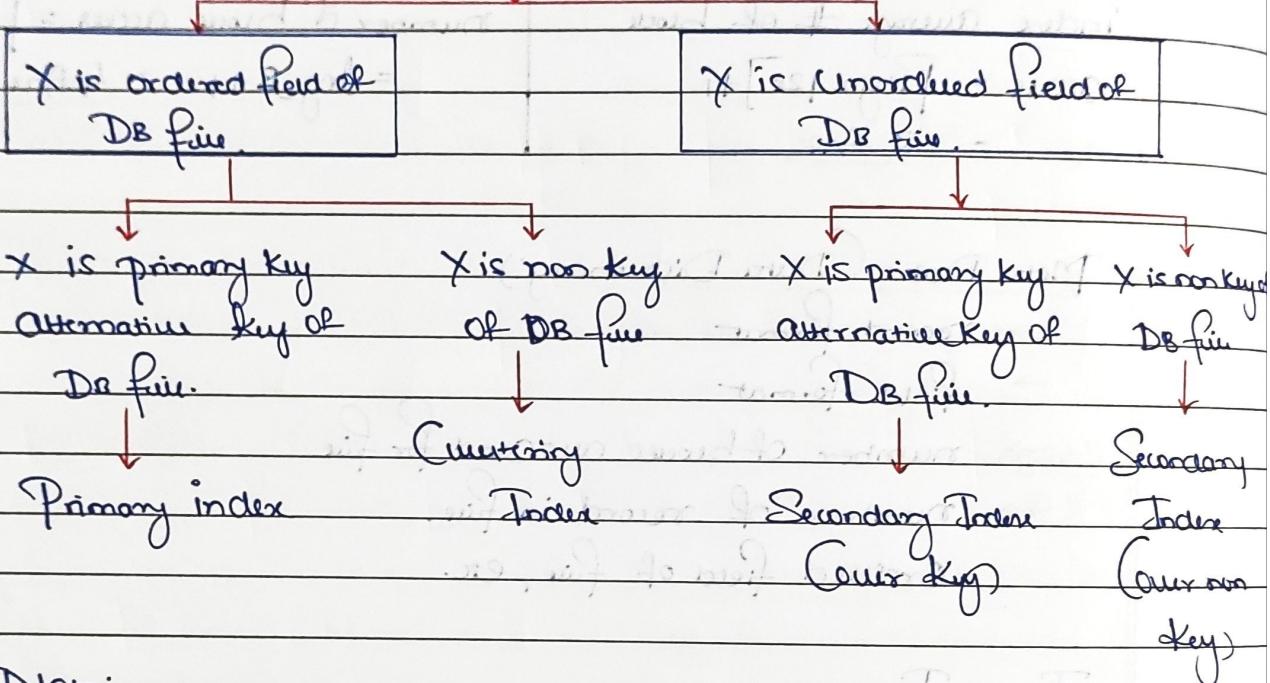
Multilevel Index

Dynamic multilevel index
 Using B and B⁺ tree.

1. Primary index (key + ordered DB file)
2. Clustering index (non key + ordered file)
3. Secondary index (non key + unordered file)

Type of index

(X : Search key)



Note:

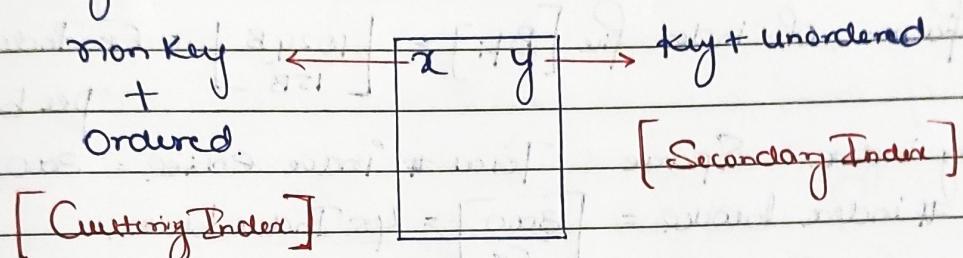
- * At most one Primary index possible per relation (Table).
- * More than one secondary index possible.
- * In a relation either CI or PI any one possible not both.

Primary Index :

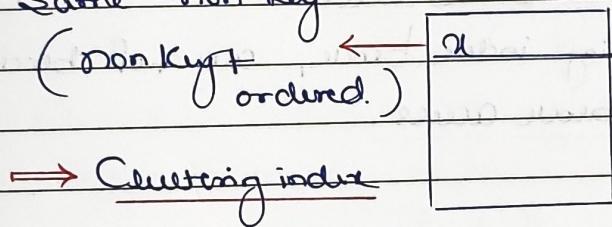
- * Ordered file with two fields:
 - * Primary key, $k(i)$
 - * Pointer to a disk block, $P(i)$
- * One index entry in the index file for each block in the data file.
- * Pointer may be dense or sparse.
 - * Dense index has an index entry for every search key, value in the data file.
 - * Sparse index has entries for only some search values.

- Access cost to access record using P_1 with multilevel index: (K_{11}) blocks.
- Primary index can be Dense or Sparse (Sparse is mostly used).
- For any database relation atmost one P_2 is possible (because of index over ordered field).

Q1. Data records based on non key (x) and index built over key field (y) of DB Table:



Q2. Data records ordered based on non key and index build over same - non key.



Q3. Suppose that we have ordered file of 30,000 records, stored on a disk with block size 1024 Byte (Record Size) file records are of fixed length and unspaced of size 100 Byte (Record Size) and suppose that we have created a primary index on the key field of the file of size 9 Byte and Block Pointer of size 6 Byte then find the average number of Block Access to search for a record using with and without index?

Ans. Records = 30,000 Block Size = 1024 B, Record Size = 100 B,
 Key = 9 B, Pointer = 6 Byte, [Ordered file, fixed length, primary index]
 Without-Index

$$\text{Block Factor of DB file} = \left\lceil \frac{1024 \text{ B}}{100 \text{ B}} \right\rceil = 10 \text{ Record/block}$$

Total number of DB record = 30,000

Total number of data block [B] = $\left\lceil \frac{30,000}{10} \right\rceil = 3000$ Data Blocks

Ordered file: To access a record avg # of

block access = $\lceil \log_2 B \rceil = \lceil \log_2 3000 \rceil = 12$ Block access

With Index

One index record size = 9+16 = 15 Byte

Block factor of index file [BF_i] = $\left\lceil \frac{1024B}{15B} \right\rceil = 68$ Index entries per block

Primary index: Sparse : Total # index entries = 3000 (# data blocks)

Total # index blocks = $\lceil \frac{3000}{68} \rceil = 45$ Index blocks.

To access index block avg # of block access = $\lceil \log_2 B_i \rceil$
 $= \lceil \log_2 45 \rceil = 6$ Block access

To access ix record using index block, avg # of block access
 $= \lceil \log_2 45 \rceil + 1 = 7$ block access

Clustering Index

Clustering Field

- * File records are physically ordered on a non key field assuming a distinct value for each record.

Ordered file using two fields

- * Same type as Clustering field

- * Disk block pointer

- * Clustering index mostly Sparse index (Dense CI also possible if each cluster contains one record).

- * Atmost One CI is possible for any database relations (Ordering required)

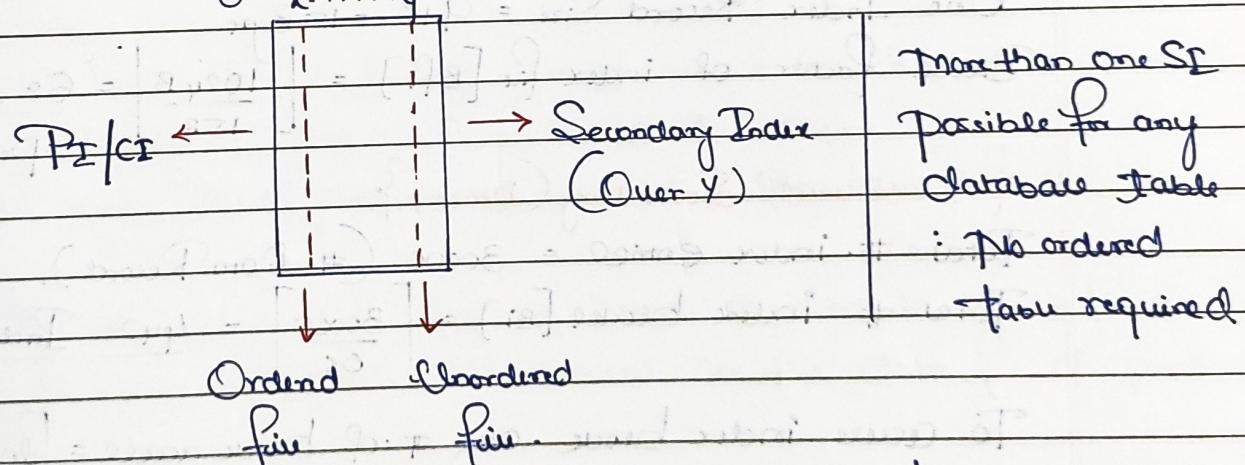
- * For any DB relations can build either PI or CI but not both.

Secondary Indexes

- * A Secondary index provides a secondary means of accessing a file for which some primary access already exists.
- * The Secondary index may be on a field which is a Candidate key and have a unique value in every record, or a non-key with duplicate values.
- * The index is an ordered file with two fields:
 - ① The first field is of the same type (datatype) as some non ordering field of that data file is an indexing field.
 - ② The Secondary field is either a block pointer or a record pointer. There may be many Secondary indexes (and hence indexing fields) for the same file.
- * Includes one entry for each record in the data file, hence it is a dense index.

→ Search Key: Unordered field and key/non-key

→ Secondary Possible way to access data using index even PI/CI indexes already exists.



- I/O cost to access record by using SI over key using PIAI is $(K+1)$ blocks
- I/O cost to access record of some non-key using SI over non-key using PI: { $K + \#$ of blocks of DB equal to # of pointers in given redirect page }.

Q4.

Consider a Secondary index on the key field of the file of question number 1, then find the average number of block access to access a record using with and without index?

Ans. # Records = 30,000, Block size = 1024 B, Record Size = 100 B

Key = 9 B, Pointer = 6 B.

{ unordered file, Secondary (Dense Index) }

Without Index.

$$\text{Block factor of DB file} = \left\lceil \frac{1024 \text{ B}}{100 \text{ B}} \right\rceil = 10 \text{ Record/block}$$

(Unspanned)

$$\text{Total number of DB record} = 30,000$$

$$\text{Total number of Data Block [B]} = \left\lceil \frac{30000}{10} \right\rceil = 3000 \text{ Data Block}$$

Ordered File

$$\text{To access a record avg \# block access} = \frac{B}{2} = \frac{3000}{2} = 1500 \text{ Block access}$$

With Index

$$\text{One Index Record Size} = 9+6 = 15 \text{ Byte}$$

$$\text{Block factor of index file [Bfi]} = \left\lceil \frac{1024 \text{ B}}{15 \text{ B}} \right\rceil = 68 \text{ Index Entries Per block.}$$

{ Unordered, Secondary (Dense) }

$$\text{Total \# index entries} = 30000 \text{ (\# Data Record).}$$

$$\text{Total \# index blocks [Bi]} = \left\lceil \frac{30000}{68} \right\rceil = 442 \text{ Index blocks.}$$

$$\begin{aligned} \text{To access index block avg \# of block access} &= \lceil \log_2 B_i \rceil \\ &= \lceil \log_2 442 \rceil = 9 \text{ Block access} \end{aligned}$$

$$\begin{aligned} \text{To access record using index avg \# block access} &= \lceil \log_2 B_i \rceil + 1 \\ &= \lceil \log_2 442 \rceil + 1 = 10 \text{ Block access} \end{aligned}$$

Multilevel Index

- * If index does not fit in memory, access becomes expensive
- * Solution: Treat index kept on index as a sequential file and construct a sparse index on it.
 - Outer index: a sparse index of the basic index
 - Inner index: the basic index file.
- * If even outer index is too large to fit into main memory, yet another level of index can be created and so on.
- * Indices at all levels must be updated on insertion or deletion from the file.
- Multilevel indices are designed to greatly reduce remaining search space as search is conducted.
- Index file: Considered first (or base level) of a multilevel index
- Second level: Primary index to the first level
- Third level: Primary index to the second level.

Note: * We can repeat the above process until index entries fit into one block.

* If there are n levels in multilevel index then the number of block access to search for a record = $n+1$
 (At each level one index block + 1 data block).

Qs. Find the average number of block access required to search for a record if multi-level index is created on the data file of previous question.

Ans. Block factor of Index file = 68 entries per block

I Level: Total no. of index blocks = 442 index blocks

II Level: No. of index records (entries) = 442 (Sparse no. of 1st level blocks) and Block factor = 68 entries per block.

Total number index blocks = $\lceil \frac{442}{68} \rceil = 7$ index blocks.

III Level: Number of index record = 7 (number of 2nd level blocks)
 Total number of index block = $\lceil \frac{7}{68} \rceil = 1$

$$\text{Average no. of blocks access} = 1+1+1+1 = 4$$

Q6

A clustering index is defined on the file, which are of type:

Ans

Non Key and Ordering:

Q7.

Consider a file A file of 16384 records. Each record is 32 byte long and its key field is of size 6 byte. The file is ordered on a non key field, and the file organization is unspaced. The file is stored in a file system using block size 1024 byte, & size of block pointer is 10 byte. If the secondary index is built on the key field of the file, and a multi-level index scheme is used to store the secondary index, the number of first-level and second-level block in the main level index are respectively are:

Ans.

#. Records = 16,384 [2¹⁴] Record size = 32 B, Key = 6 B, Bp = 10 B,

Block size = 1024 B, Secondary Index (Dense) & multi-level index:

One index record size = 6 + 10 = 16 Byte

$$\text{Block Factor of Index File } [BFI] = \left\lceil \frac{1024 B}{16 B} \right\rceil = \frac{2^{10}}{2^4} = 2^6 = 64 \text{ Index entries per block}$$

Secondary Index: (Dense).

Total #. index entries = 16,384 (# DB records)

$$\text{I Level: Total # index block} = 2^4(16,384) = 2^8 = 256 \text{ Index Block}$$

$$2^6(64)$$

II Level: #. index Entries = 256 (# Ist level blocks)

BFI = 64 Index entries per block

$$\text{Total # index blocks} = \left\lceil \frac{256}{64} \right\rceil = 4$$

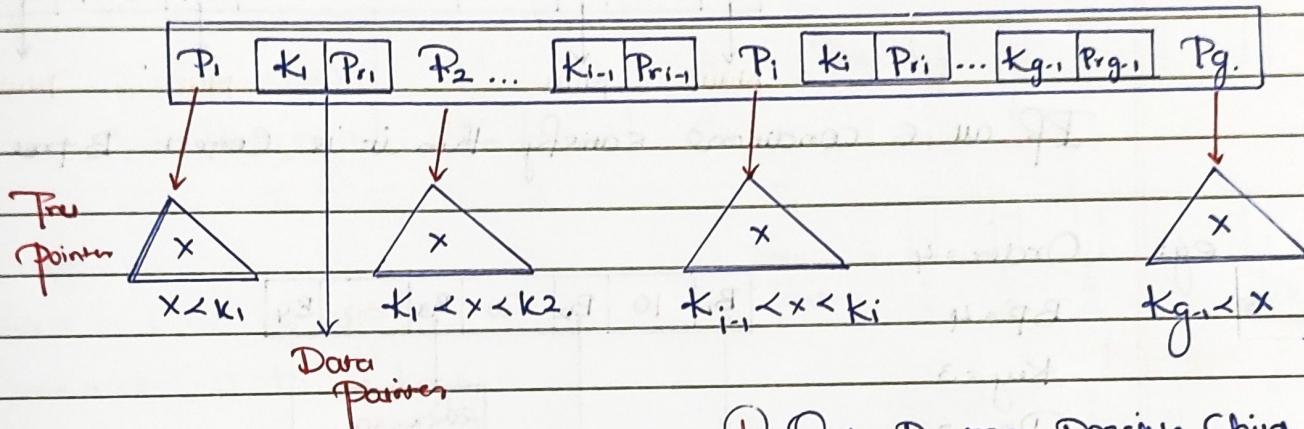
$$\therefore \text{Ans} = \{256, 4\}$$

III Level: Total # index entries = 1, & $B_F = 64$ Index Entries per block.

$$\text{Total # index blocks} = \left\lceil \frac{4}{64} \right\rceil = 1$$

\therefore So total 3 level required.

B Tree Structure:



① Order P : max possible Child

Pointers [degree] can
Store in B Tree node

P: Block/Child/Index Pointer

Pr: Read/Record/Data Pointer

k: Search Key

B-Tree Definition

- ② Every internal node except the root node Contains atleast (min) $\lceil P/2 \rceil$ block pointer (min keys $\lceil P/2 \rceil - 1$) And maximum P block pointer (maximum keys $P-1$).
- ③ Root Can Contain atleast 2 block pointer (min 1 key in root node) and maximum P block pointer (max $P-1$ keys).
- ④ Keys within the node Should be in ascending order.
- ⑤ Each leaf node Should be at same level.

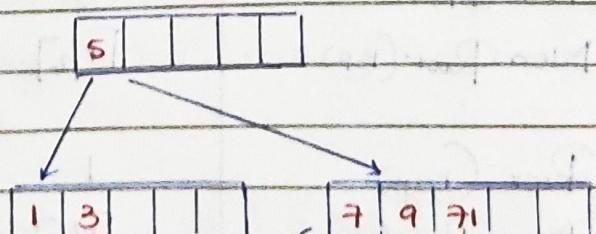
eg. Order = 6

$$\rightarrow 1, 5, 3, 7, 9, 11, 13, 16 \rightarrow$$

1	3	5	7	9
---	---	---	---	---

$$\text{Min Key} = \lceil 6/2 \rceil - 1 = 2$$

$$\text{Maximum keys} = 6 - 1 = 5.$$



Order = P

BP = P

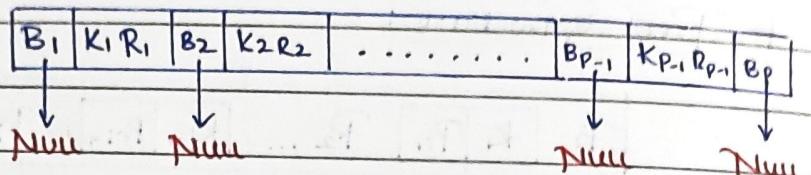
key = $P-1$

RP = $P-1$

any sick can have
more keys.

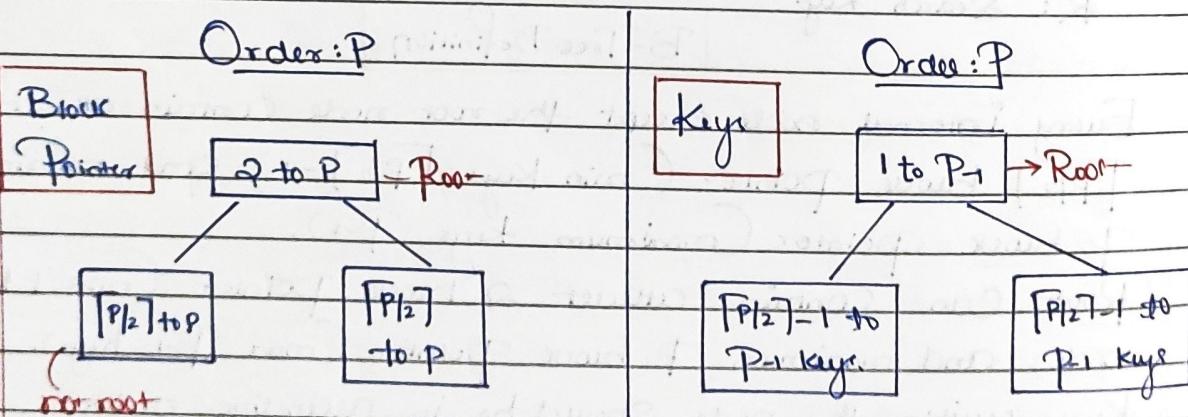
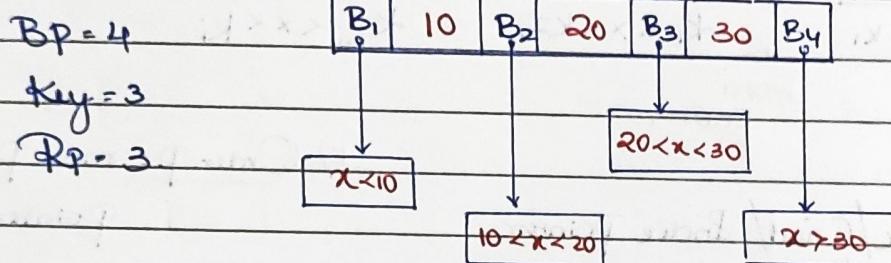
- * Order : P
- * $\text{Max Keys} = P - 1$
- * $\text{Min Keys} = \lceil P/2 \rceil - 1$
- * $\text{Max BP} = P$
- * $\text{Min. BP} = \lceil P/2 \rceil$

(6) Structure of leaf node:



If all 6 conditions satisfy then it is called Complete B-tree of order P .

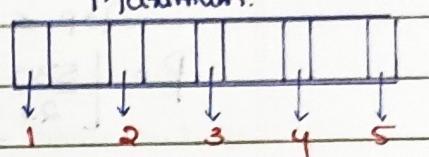
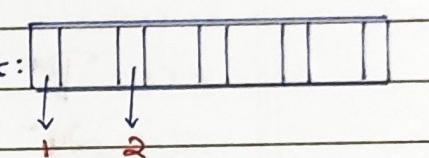
e.g: Order = 4



Order : P	Minimum	Maximum
Root (BP):	2	P
Non Root (BP):	$\lceil P/2 \rceil$	P
Root (keys)	1	$P-1$
Non Root (keys)	$\lceil P/2 \rceil - 1$	$P-1$

Note: Root cannot have only one block pointer.

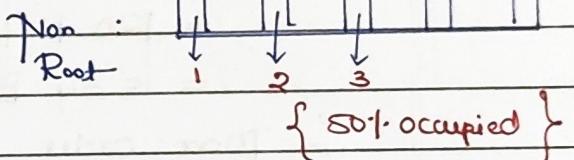
Order: P [5].

	Maximum.	Minimum.	Maximum.
Root (Bp)	2	5	
Non Root (Bp)	$\lceil \frac{5}{2} \rceil = 3$	5	
Root (Key)	1.	4	Root: 
Non Root (key)	$\lceil \frac{5}{2} \rceil - 1 = 2$	4.	

Order: 5

Root: 1 to 4 keys

Non root: 2 to 4 keys

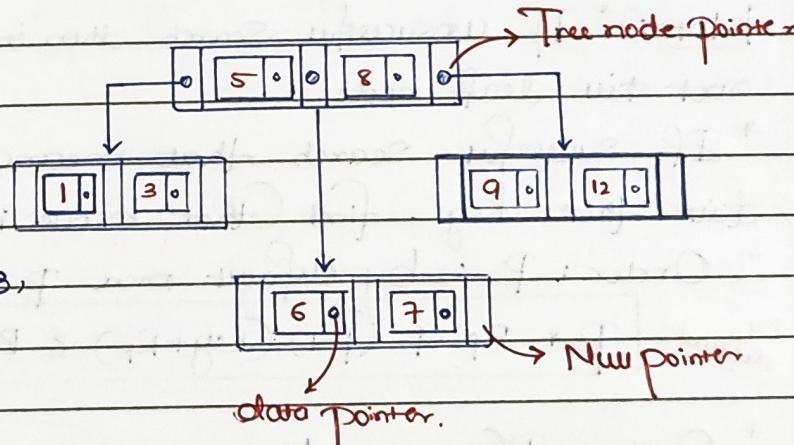


→ A B-Tree of order
 $P=3$. The values

were inserted in

the order 8, 5, 1, 7, 3,

12, 9, 6.



Q1. Consider a table in a relational database with a key field K. A B-tree of order p is used as an access structure on K, where p denotes the maximum number of tree pointers in B-tree index node. Assume that K is 10-byte long, disk block size is 512 bytes. Each data pointer size P_d is 8 bytes long and each block pointer P_b is 5 bytes long. In order for each B-tree node to fit into a single disk block, the maximum value of p is

Ans $K = 10 \text{ Byte}, P_d = 8 \text{ B}, P_b = 5 \text{ B}, \text{Block Size} = 512 \text{ Byte}$

Order: P.

$$P \geq B_p + (P-1)[K_{\text{bytes}} + P_d] = \text{Block Size.}$$

$$\begin{aligned}
 P^* S + (P-1)[10+8] &\leq 512 \\
 5P + (P-1)[18] &\leq 512 \\
 5P + 18P - 18 &\leq 512 \\
 23P &\leq 530 \\
 P = \left\lfloor \frac{530}{23} \right\rfloor &= \lfloor 23.04 \rfloor = \underline{\underline{23}}
 \end{aligned}$$

If we take 24 then we get the values of order ($P=24$) - Then it exceeds by block size 512 Byte
 Or
 It cannot accommodate in a block of size 512 Byte.

$$\begin{aligned}
 P^* Bp + (P-1)(Key + Rp) \\
 24^* Bp + (24-1)(Key + Rp) \\
 24^* S + 23(18) \\
 = 120 + 414 \\
 = 534 \text{ Byte, But given block size is 512 Byte.} \\
 \therefore \text{More order is } \underline{\underline{23}} \text{ not 24.}
 \end{aligned}
 \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \rightarrow \text{Block Size (512 Byte)}$$

Note: * If unsuccessful search then it means we traversed from root till leaf node.

* If successful search that means we traverse from root to till find key. find (that block in which key is present).

* Order: P : by default max P block pointer

By default $P^* Bp + (P-1)(Key + Rp) \leq \text{Block size.}$

* Order: P(key), key order: P

$$(P+1) Bp + P(Key + Rp) \leq \text{Block size.}$$

Order: P	Height/level	Min # nodes	Min # Bp	Min # Keys	* Level = height + 1
Minimum	0/1.	1.	2	1.	
	1/2.	2.	2 $\lceil P/2 \rceil$	$2(\lceil P/2 \rceil - 1)$	
	2/3.	$2 \lceil P/2 \rceil$	$2 \lceil P/2 \rceil^2$	$2 \lceil P/2 \rceil (\lceil P/2 \rceil - 1)$	
	:	:	:	:	
	b/b+1	$2 \lceil P/2 \rceil^{b-1}$	$2 \lceil P/2 \rceil^b$	$2 \lceil P/2 \rceil (\lceil P/2 \rceil - 1)$	

→ Total Minimum = $\left[1 + 2(\lceil \frac{P}{2} \rceil - 1) + 2\left[\frac{P}{2}\right](\lceil \frac{P}{2} \rceil - 1) + \dots + 2\left[\frac{P}{2}\right]^{h-1} (\lceil \frac{P}{2} \rceil - 1) \right] \rightarrow GP \text{ Formula.}$

Order: P

Maximum.

Height/Level.	Max #. Nodes	Max #. BP	Maximum Keys
0/1	1	P	(P-1)
1/2	P	P^2	$P(P-1)$
2/3	P^2	P^3	$P^2(P-1)$
3/4	P^3	P^4	$P^3(P-1)$
⋮	⋮	⋮	⋮
b/b+1	P^b	P^{b+1}	$P^b(P-1)$

Total Maximum = $(P-1) + P(P-1) + P^2(P-1) + P^3(P-1) \dots + P^b(P-1)$
Keys = $(P-1)(1 + P + P^2 + P^3 \dots + P^b)$.
⇒ Apply GP Series

Q1. Order: 5
Level: 5
= $(5-1) + (5 \times 4) + (5^2 \times 4) + (5^3 \times 4) + (5^4 \times 4)$
= $4 + 20 + 100 + 500 + 2500 \Rightarrow 3124$

Total max no. of keys = 3124.

Q2. Order: 23 : Consider a BTree 69% Full, then minimum number of keys?

No. of BP per node = 23×0.69
= 16

Fanout = 16

(P) Order: 16

(P-1) Keys = 15 $(\lceil \frac{P}{2} \rceil - 1) \quad \lceil \frac{P}{2} \rceil$

Note: Maximum level: At each level min # key & min # BP

Minimum level: At each level maximum # key & max # BP
{ Applicable for both B and B1 Tree }

Advantage of B-Tree:

- ① B Tree index best Suitable for random access of some records.
Select *
From R
Where $A = 24$:

One record
access.

$$\begin{aligned} \text{I/O cost} &: K+1 \text{ blocks} \\ &= [\log_P n] + 1 = \Theta(\log_P n) \end{aligned}$$

Disadvantage of B-Tree:

- ① B Tree index not best Suitable for Sequence access of range of records.

SELECT * FROM R

WHERE $A \geq 30$ and $A \leq 85$;

I/O cost: $\alpha(\log_P n + 1) + \text{cost of unsuccessful}$

more access cost

(Unordered from DB).

Range of records

B+ Trees

- All keys are available at leaf node
- Internal node: no read pointer
- Leaf node contains key & $R_p + 1$ Block pointers.

Difference between BTree and B+ Tree.

- ① BTree Internal Node:

B_1	K_1	R_1	B_2	K_2	R_2	B_n	K_p	R_p	B_p
-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------

Order: P (maximum #. BP)

$B_p = P$, $Keys = P-1$, $R_p = P-1$.

$$P + BP + (P-1)(Key + R_p) \leq \text{Block size}$$

- B+ Tree Internal Node:

B_1	K_1	B_2	K_2	B_p	K_p	B_p
-------	-------	-------	-------	-------	-------	-------	-------

Order: P (maximum #. pointers)

$P * BP + (P-1) \text{key} \leq \text{Block size}$

- ② BTree Leaf Node:

B_1	K_1	R_1	B_2	K_2	R_2	B_p	K_p	R_p	B_p
-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------

- B+ Tree Leaf Node:

K_1	R_1	K_2	R_2	K_p	R_p	•
-------	-------	-------	-------	-------	-------	-------	---

1 Block
pointer

Order : P (max #. of BP)

$Bp = P$, $key_1 = P-1$, $Rp = P-1$.

$$P * Bp + (P-1)(Key + Rp) \leq \text{Block Size}$$

Order : P (Maximum # pointers)

$$(P-1)(key + Rp) + 1Bp \leq \frac{\text{Block Size}}{1024}$$

- Note:
- * In a B Tree internal node and leaf node have same structure (order of internal node is same as order of leaf node in B tree)
 - * In a B^+ Tree internal node and leaf node have different structure (order of internal node is not same as order of leaf node).

eg:1 Block size = 1024 Byte, Key = 6 Byte, $Bp = 8B$, $Rp = 10$ Byte order of internal node and leaf node in a B Tree?

$$P * Bp + (P-1)(Key + Rp) \leq \text{Block Size}$$

$$P * 8 + (P-1)(6+10) \leq 1024$$

$$8P + 16P - 16 \leq 1024$$

$$P * 24 \leq 1024 \Rightarrow \left\lfloor \frac{1024}{24} \right\rfloor = 43.$$

Order

Internal node : 43

Leaf node : 43

$$\therefore \underline{P = 43}$$

eg:2 Consider block size = 1024 Byte, Key = 6 Byte, $Bp = 8B$, $Rp = 10$ Byte
Order of internal and leaf node in a B^+ Tree.

Order of internal node

$$P * Bp + (P-1) Key \leq \text{Block Size}$$

$$P * 8 + (P-1) 6 \leq 1024$$

$$8P + 6P - 6 \leq 1024$$

$$14P \leq 1030$$

$$P = \left\lfloor \frac{1030}{14} \right\rfloor = \underline{73.}$$

Order of leaf node in B^+ Tree

$$(P-1)(key + Rp) + 1Bp \leq \text{Block Size}$$

$$(P-1)(6+10) + 8 \leq 1024$$

$$16P - 16 + 8 \leq 1024$$

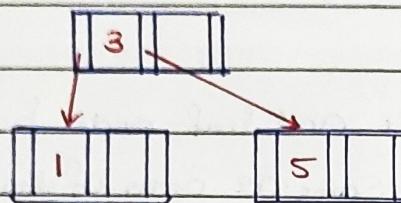
$$16P \leq 1024$$

$$P = \left\lfloor \frac{1024}{16} \right\rfloor = \underline{64}$$

eg: Keys = 1, 3, 5

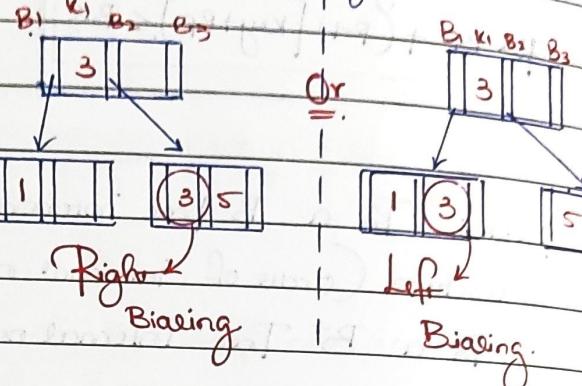
B Tree

Order : 3 Max keys = 2.



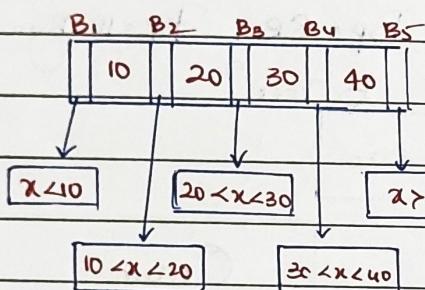
B^+ Tree

Order : 3, Max keys = 3.

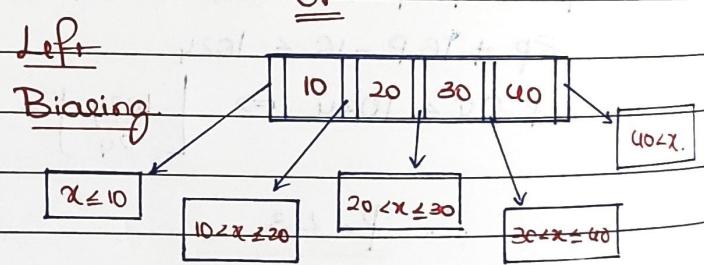
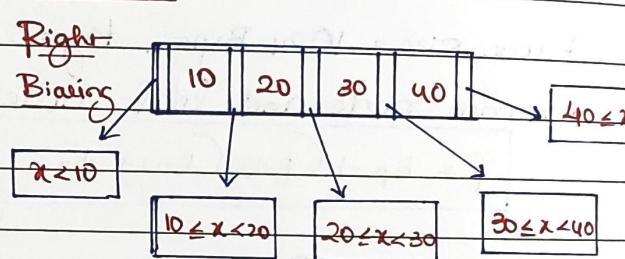


eg.

B Tree: Order : 5



B^+ Tree Order : 5



B^+ Tree Definition

(1) Structure of Internal node:



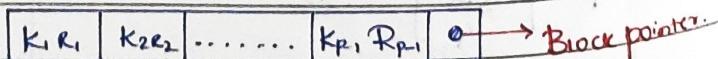
Left biasing $\Rightarrow x \leq K_1 \quad K_1 < x < K_2$

or

Right biasing $\Rightarrow x < K_1 \quad K_1 \leq x < K_2$

$P * B_p + (P-1) \text{ key} \leftarrow \text{Block Size}$

(2) Structure of leaf node:



$(P-1) [\text{key} + R_p] + 1 B_p \leftarrow \text{Block Size}$

Remaining 4 points are same as B^+ tree

Q1. Why internal node not having Rp?

Ans. As all keys are available at leaf node so read pointer not required in internal node.

Q2. Why leaf node have only key and Rp only?

Ans. Because all the keys are available in leaf node according to B⁺ definitions so to read those keys (data) read/data pointer is required.

Q3. Why leaf node having only 1 block pointer?

Ans. B⁺ Tree is suitable for range query, and in B⁺ tree all keys are available in the leaf node and if 1 block connect to next node at leaf (like linked list) so I/O cost is reduced while performing range query.

Q4. When all keys are available at leaf node, then why some keys are available in non leaf node?

Ans. Yes there are some keys available in the non leaf node but they are very useful in searching purpose. It is only used for searching purpose, they don't have Read (Record/data) pointer in non leaf node.

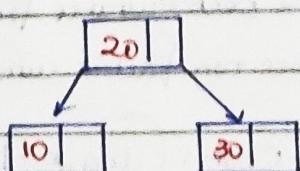
Creation/Insertion in B⁺ Tree

① In B⁺ Tree insertion start from leaf node.

② Split of internal node is same as split of internal node in B Tree.

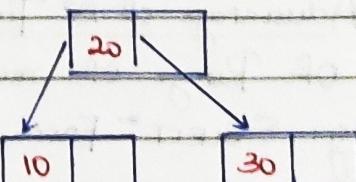
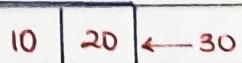
Internal Node

B Tree:



Internal Node

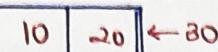
B⁺ Tree:



3. Split of leaf in B^+ Tree is different from B Tree.

Leaf Node:

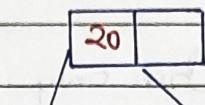
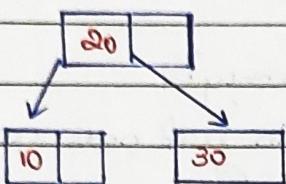
B Tree.



Leaf Node:



B^+ Tree



Left Biasing.

Right Biasing

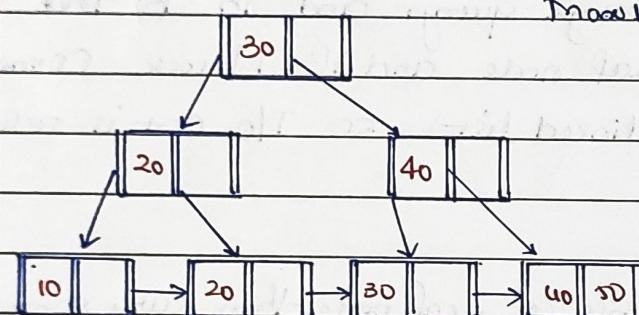
eg:

Keys: 10, 20, 30, 40, 50.

Order: 3

{ Right Biasing }

Max Key = 2



Important Points:

- * Every key of leaf node is not present in non-leaf node.
- * example: Keys = 10, 20, 30, 40, 50 in leaf node but all keys are not present in non-leaf node (except 10, 50 all are present in previous example).
- * Some key of leaf node is present in non-leaf node.
- * But every key of non-leaf node must be present in leaf node (i.e. in B^+ Tree all keys are available at the leaf node).
- * To B^+ Tree: Leaf node key order:

$$P(\text{Key} + R_p) + 1 B_p \leq \text{Block Size}$$

Advantages of B^+ Tree: B^+ Tree index better suitable for sequence access of Range of records. (Range Queries run faster with B^+ Tree index).

e.g. SELECT * From R

WHERE A \geq 30 and A $<$ 85;

\rightarrow I/O Cost: $(\log p^n - 1) z + x$ for index
for DB file.