

Relation Algebra

Basic operators:

Projection (π) , Selection (σ) , Cross-product (\times)
Union (\cup) , Rename (ρ) , Set-difference (-)

Derived operators:

Join (\bowtie) , Intersection (\cap) , Division ($/, \div$)

Projection (π): used to select columns.

Selection (σ): used to select rows.

Cross-product (\times): More combinations of each row in A
with every row in B.

$$\text{Degree}(A \times B) = D(A) + D(B)$$

$$\text{Cardinality}(A \times B) = C(A) \times C(B)$$

Set-difference (-): $A - B \Rightarrow$ In A but not in B

$$A - B \neq B - A, \text{ (It is not commutative).}$$

Union (\cup): ① No. of columns must be same.

② Domain of every column must be same.

$$\text{Degree}(A \cup B) = D(A) - D(B)$$

$$\text{Cardinality}(A \cup B) = C(A) + C(B)$$

Division: $A(x, y) / B(y) =$ It results x values for that
there should be a tuple $\langle x, y \rangle$ for every Y value
of relation B. e.g. Find Sid of students who are enrolled in all courses.

$$\pi_{\text{sid}}(\text{Enrolled}) - \left[\pi_{\text{sid}} \left((\pi_{\text{sid}}(\text{Enrolled}) \times \pi_{\text{cid}}(\text{course})) - \text{Enrolled} \right) \right]$$

Indexing

Indexing is used to reduce I/O cost i.e. reducing no. of blocks moved to RAM.

Q. Block size = 1000 Bytes (without indexing)

Record size = 250 Bytes

Total records = 10,000 (unordered)

Avg. time complexity to search = ?

No. of records we can put in one block =

$$\frac{1000 \text{ (BS)}}{250 \text{ (RS)}} = 4 \text{ records}$$

$$\text{No. of blocks required} = \frac{10,000}{4} = 2500$$

I/O cost; Best case = 1

worst case = 2500

$$\text{Avg. case} = \frac{2500}{2} = 1250$$

In case of ordered data binary search can be used, then time complexity = $\log_2 2500 = 12$

Q what is avg time complexity to search a record from index table if index table entry = 20 B(key + pointer)
10B 10B

Block size = 1000 Bytes

Index table entry = 20 Bytes

No. of records that we can put in one block of index table = $\frac{1000}{20} = 50$ records.

$$\text{No. of index blocks required (worst)} = \frac{2500}{50} = 50$$

$$\text{Time complexity (I/O)} = 50[\log_2 50] = 6 + \frac{1}{1} = 7$$

$$\text{If Dense then time, } = \frac{10,000}{50} = 200 \text{ records read actual block}$$

$$\text{Time} = [\log_2 200] + 1 = 8 + 1 = 9$$

unordered data → Dense index → record for each entry
 ordered data → Sparse index → record for each block.

Types of Indexes

		key	non-key	
1) Primary				
2) clustered	ordered	Primary Index	clustered Index	
3) Secondary	file	(sparse)	(sparse)	At most 1
	unordered file	secondary Index	secondary Index	

Primary Index

- Key should be primary key and file must be ordered.
- Sparse index is generally used i.e one pointer for each block.
- First record of each block is called anchor record.
- No. of entries in index table = no. of blocks in hard disk.
- Search time = $\log_2 N + 1$ where N is no. of blocks in index table.

Clustered Index

- Non key but data is sorted.
- one entry in index table for each unique value.
- sparse
- Block pointer is used in a case when a repeating value spans more than one block.
- Time = $(\log_2 N + 1) \rightarrow$ minimum, may exceed because of block pointer.

Secondary Index

-) When data is unordered (key or nonkey)
-) In a case when column is key (unique) but not ordered, Dense index is created which is sorted, so binary search can be used.
-) No. of records in index = no. of records in disk.
-) Time $\log_2 N + 1$ where N = no. of record blocks in indexable.

-) In a case when column is neither key nor ordered, index is created for each unique value and is sorted.
-) In addition a 'block of record pointers' is maintained, which points to all duplicate value records. (Intermediate layer).
-) Generally dense because of Intermediate layer.
-) Min. time = $(\log_2 N + 1 + 1)$.

Normalization

It is a technique to ~~redundancy~~ remove or reduce redundancy from a table.

Insertion anomaly, deletion anomaly and updation anomaly.

First normal form:

Table should not contain a multivalued attribute.

Closure method: used to find all primary keys of the table.

$R(A B C D E)$

$FD = \{A \rightarrow B, BC \rightarrow D, E \rightarrow C, D \rightarrow A\}$

$E^+ = EC$ (E wasn't present on right hand side)

$AE^+ = AEBCD$

$CK = \{AE, DE, BE\}$

$DE^+ = DEACB$

Prime att. = A E D B

$BE^+ = BECDA$

Non-Prime att. = C

$CB^+ = CE \times$

Functional Dependency:

$x \rightarrow y$

x determines y.

For same value of determining attribute, value of determined attribute should be same

Reflexivity: If Y is a subset of X then $X \rightarrow Y$ (trivial)

Augmentation: If $X \rightarrow Y$, then $XZ \rightarrow YZ$

Transitive: If $X \rightarrow Y$ and $Y \rightarrow Z$ then $X \rightarrow Z$

Union: If $X \rightarrow Y$ and $X \rightarrow Z$ then $X \rightarrow YZ$

Decomposition: If $X \rightarrow YZ$ then $X \rightarrow Y$ and $X \rightarrow Z$

Pseudotransitivity: If $X \rightarrow Y$ and $WY \rightarrow Z$ then $WX \rightarrow Z$

Composition: If $X \rightarrow Y$ and $Z \rightarrow W$ then $XZ \rightarrow YW$

Second normal form

- 1) Table or relation must be in 1st normal form.
- 2) There should be no partial dependency i.e all non prime attribute of a table should be fully functional dependent on candidate key.

$R(ABCDEF)$

FD { $C \rightarrow F$, $E \rightarrow A$, $EC \rightarrow D$, $A \rightarrow B$ }

$EC^+ = ECDFAB$

CK = { EC }

Prime attribute { E, C }

Non-prime attribute { A, B, D, F }

Here F is determined by C i.e a part of primary key.

Hence, this table is not in 2nd normal form.

Partial dependency: L.H.S should be a proper subset of candidate key and R.H.S should be a non prime attribute.

third normal form

- 1) The relation must be in second normal form
- 2) There should be no transitive dependency i.e a non prime attribute should not be determined by a non prime attribute.
- 3) For each FD, L.H.S must be a candidate key or super key or R.H.S is a prime attribute.

$R(ABCD)$

FD { $AB \rightarrow CD$, $D \rightarrow A$ }

CK $AB^+ = ABCD$

Since AB is candidate key in ($AB \rightarrow CD$) and A is a prime attribute in ($D \rightarrow A$), the relation is in third normal form.

BCNF (Boyce Codd Normal Form)

L.H.S of each FD should be a candidate key or a super key.

- third normal form always ensures 'dependency preserving decomposition' but not in BCNF.
- Both third and BCNF ensures lossless decomposition.

$R(ABCD)$

FD: $\{AB \rightarrow CD, D \rightarrow A\}$

$AB^+ = ABCD$

$DB^+ = DBAC$

CK = {AB, BD}

PA = {A, B, D}

NPA = {C}

ABCD

DA BCD

{ $D \rightarrow A\}$ }

$BD^+ = BDAC$

{ $BD \rightarrow C\}$

Here the relation is in third normal form but not in BCNF because in ($D \rightarrow A$) 'D' is not candidate key. Hence we take out D^+ i.e. DA from relation.

Also $D \rightarrow A$ is preserved but $AB \rightarrow CD$ is not. Hence we have to make new dependency $BD \rightarrow C$.

In lossy decomposition, spurious tuples gets generated. To prevent that the common attribute should be a prime candidate key or super key of either R_1 or R_2 or both.

If A is the candidate key: If in $R(ABC)$, then $R_1(AB), R_2(AC)$

1) $R_1 \cup R_2 = R$

3) $R_1 CK$ or $R_2 CK$ or both.

$AB \cup AC = ABC$

$ABC = ACC$

a) $R_1 \cap R_2 = \emptyset$

$AB \cap AC = \emptyset$

$A \neq \emptyset$

Equivalence of functional dependency

$$Y \text{ } X \text{ covers } Y \quad X \supseteq Y$$

$\Leftrightarrow Y \text{ } X \text{ covers } X \quad Y \supseteq X \subseteq Y$

$\boxed{X \equiv Y}$

$$X = \{A \rightarrow B, B \rightarrow C\}, \quad Y = \{A \rightarrow B, B \rightarrow C, A \rightarrow C\}$$

$Y \text{ } X \text{ covers } Y$

In X

$$A^+ = A \quad (\text{reflexive})$$

$$A^+ = AB \quad (\text{reflexive})$$

$$A^+ = ABC \quad (\text{transitive})$$

In X

$$B^+ = B \quad (\text{reflexive})$$

$$B^+ = BC \quad (\text{reflexive})$$

So, A is covering $A \rightarrow B$ & $A \rightarrow C$

Hence X -covers Y

So, B is covering $B \rightarrow C$

$Y \text{ } Y \text{ covers } X$

In Y

$$A^+ = ABC$$

So A is covering $A \rightarrow B$

Hence Y covers X.

In Y

$$B^+ = BC$$

So, B covers $B \rightarrow C$ in X

$$1 \quad X \supseteq Y$$

$$2 \quad Y \supseteq X$$

$$\Rightarrow \boxed{X \equiv Y}$$

B - Tree

→ Block pointer points to other node in tree.

→ Data / Record pointer points to actual record stored in harddisk.

→ For a tree with order p

$$\text{No. of block pointers} = p$$

$$\text{No. of max. children} = p$$

$$\text{Keys} = p - 1$$

$$\text{Record pointers} = p - 1$$

→ Keys are inserted in sorted order

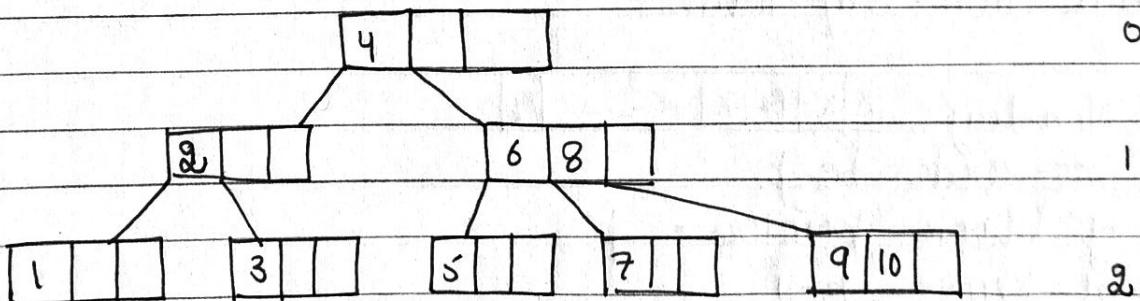
children	Root	Intermediate Node
Max	p	p
Min	2	$\lceil \frac{p}{2} \rceil$

a) Insert following keys into B-tree, if order of B-Tree = 4.

Ans. 1, 2, 3, 4, 5, 6, 7, 8, 9, 10

$$\text{Max keys} = p - 1 = 3$$

$$\text{Min keys} = \lceil \frac{p}{2} \rceil - 1 = 1$$



a) Consider a B-tree with key size = 10 bytes, block size 512 bytes, data pointers of size 8 bytes and block pointer is 5 bytes. Find the order of B-tree?

Ans. Let the order of tree be p

$$\text{no. of block pointers} = p$$

$$\text{no. of keys and data pointers} = p - 1$$

$$\text{So, } p \times B_p + (p-1) \text{key} + (p-1) R_p \leq \text{Block size}$$

$$5p + (p-1)(10+8) \leq 512 \Rightarrow 23p \leq 512 \Rightarrow p \leq 22$$

$$\text{Min: } \lceil \frac{p}{2} \rceil = 12 \quad \text{i.e. } 4 \times B_p \text{ or children so } p = 23 \quad p \leq 23.04$$

B-Tree

- 1) Data is stored in leaf as well as internal nodes.
- 2) Internal nodes have key and record pointers.
- 3) Searching is slower and deletion complex.
- 4) No redundant search key present.
- 5) Leaf nodes are not linked together.
- 6) Leaf node have Null B_p , B_{lp} .

B^+ -Tree

- 1) Data is stored only in leaf node.
- 2) Record pointers are not there due to which more keys can be stored.
- 3) Searching is faster and deletion easy (directly from leaf node).
- 4) Redundant keys may present.
- 5) Linked together like linked list.
- 6) ~~B_p are not present at leaf~~
- 7) only one B_p for linking.

Q Consider a B^+ tree with key size = 10 bytes, block size = 512 bytes, data pointers = 8 bytes and block pointer = 5 bytes. What is the order of leaf and non-leaf nodes?

For Non-Leaf $\boxed{B_p | K | B_p | K | \dots | B_p}$

Let the order be p .

No. of block pointers = p

No. of keys = $p - 1$

So, (key size) \times $p - 1$ + $p \times B_p$ size ≤ 512

$$10(p-1) + 5p \leq 512$$

$$15p \leq 522$$

$$p \leq \frac{522}{15} \quad 34.8$$

$$\text{So } p = 34$$

For Leaf $\boxed{K | R_p | K | R_p | \dots | R_p \rightarrow}$

$$n(K + R_p) + B_p \times 1 \leq \cancel{BS}$$

$$n(10 + 8) + 5 \leq 512$$

$$n \leq 28.16 \text{ or } n = 28$$

Let n = order of leaf node i.e. max no. of (key, R_p) pair