

Normalize

Given a relation & a set of functional dependencies

$$\text{Col A} \rightarrow \text{Col B}, \text{Col C}$$

$$\text{or Col A, Col B} \rightarrow \text{Col C}$$

Decompose/normalize relation w/o info loss and so that functional dependences can be enforced.

Rules:

1. No redundancy of facts
2. no cluttering facts
3. Must preserve info
4. must preserve functional dependences

Relations might have several problems?

1. Redundancy: same value for a particular column.
2. Insertion anomaly: insert only partial info leaving null value
3. Deletion anomaly: delete a row, whole row info lost, we may just want to delete part of it for
4. Update anomaly: only update some row, causing consistency problem.

We can decompose a relation into 2. However, if not done properly, we might have 1. info loss. e.g. the connecting column is not unique, might causing extra rows that weren't in original table, this is also info loss 2. dependency loss, we can't have $\text{Col A, B} \rightarrow \text{Col C}$ when Col C is in another relation

How to decompose?

Functional Dependencies:

def: Let X & Y be sets of attributes in R

Y is functionally dependent on X in R iff for each $x \in X$ there is precisely one $y \in R[Y]$.

Full Functional Dependencies: extraneous attribute: remove it w/o changing closure.

Y is Full Functional Dependencies if Y is functional dependent on X & Y is not functional dependent on any proper subset of X

prime attribute: if it is a member of some candidate key.

- We use keys to enforce full functional dependencies $X \rightarrow Y$
- In a relation, the values of the key are unique
(let X be the key, if $X \rightarrow Y$)

How well a relation lays out

4. Normal forms

Whole set of data structures is called non-first normal form

1NF



1NF: R is in 1NF iff all domain values are atomic

2NF: 1NF & every non-key attribute is fully dependent on the key

3NF: 2NF & any Non-key attribute is non transitively dependent on the key

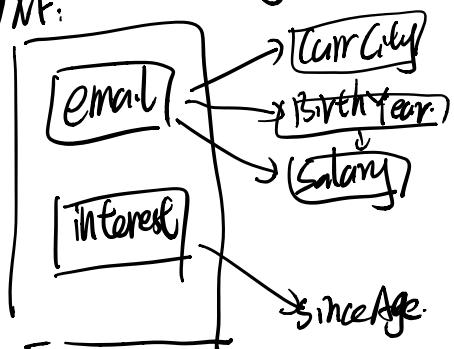
* 3CNF: every determinant is a candidate key
determinant: a set of attributes on which some other attribute is fully functionally dependent.

or 1NF: All attributes must depend on the key

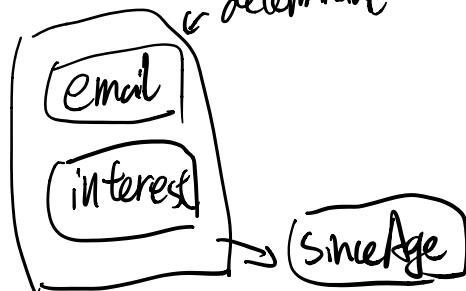
2NF: All attributes must depend on the whole key

3NF: All attributes must depend on the whole key & nothing but the key

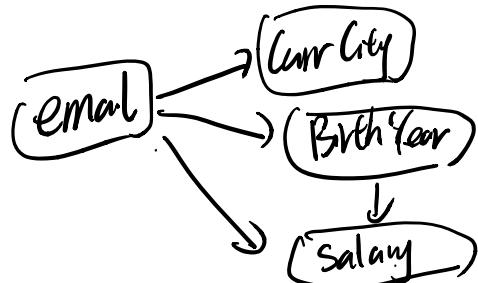
e.g. 1NF:



3CNF:



2NF:



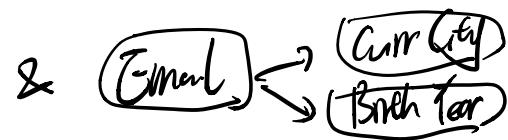
(key is email)

transitively dependent:

email \rightarrow Birth Year \rightarrow Salary

\Downarrow 3NF & BCNF





How to compute w/ Functional Dependencies:

Armstrong's rules:

reflexivity: if Y is part of X , then $X \rightarrow Y$

e.g. Email, Interest \rightarrow Interest

this is trivial.

$$\Rightarrow \varphi_{X \rightarrow YZ} \Rightarrow X \rightarrow Y$$

$$Q_2 X \rightarrow Y \quad X \rightarrow Z$$

$$\Rightarrow X \rightarrow YZ$$

$$Q_3 X \rightarrow Y, WY \rightarrow Z$$

$$\Rightarrow WX \rightarrow Z$$

augmentation: if $X \rightarrow Y$ then $W, X \rightarrow WY$

if $abA \rightarrow cdB$ then

$$abA, c \rightarrow abB, c$$

transitivity: if $X \rightarrow Y$ & $Y \rightarrow Z$ then $X \rightarrow Z$

F^+ , extends to all FD implied by F . X^+ , all attributes can be determined by X .

How to guarantee lossless Joins

$$F^+ = E^+ ?$$

$F^+ = E^+$? $F^+ = E$ equivalent

Join field must be a key in at least one of the relations

How to guarantee preservation of Functional Dependency

the meaning implied by the remaining functional dependencies must be the same

In theory, there is relation that is 3NF but not BCNF, while (and can only be decomposed to 3NF)

being lossless & dependency preserving.

It can only happen when the relation has overlapping keys

it really never happens in practice though ---

minimal FD set: 1. every dependency has 1 attribute in right hand side. 2. cannot replace any $X \rightarrow A$ w/ $Y \rightarrow A$ where $Y \not\subseteq X$ & still have a set of dependencies equivalent to P. 3. cannot remove any dependencies and still stay equivalent

Efficiency

RAM: Main memory, volatile, fast small & expensive

Access time

3×10^{-7} sec
(30 nanosec)

Disk: second memory, permanent slow big & cheap

10^{-2} sec.
(1-10 millisecond)

- CPU application can only query & update data in main memory
 - Data must be written to secondary memory after updated.
- Only a tiny fraction of a real database in main memory

We only need to consider I/O time \approx CPU time
 \approx small 10^{-7} sec.

Records: Regular User(

e.g. Email varchar (50) ← 159 bytes
sex char (1) , 8 bytes first few → end ...
Birthdate datetime
currCity varchar (50)
HomeTown varchar (50)); where record start etc.

Assume block size 4k. W/o metadata & 8% filled. \Rightarrow 320 bytes

\Rightarrow ~~20 records~~ for each block

~~blocking factor~~
Spanned representation: record broken up int 2 blocks
usually don't do that

blocks have pointer, connect to next block of this file

Transport data between disk & main memory.

Seek time: 3-8 ms Transfer time: 0.5-1 ms

rotational delay: 2-3 ms \Rightarrow Total 5-12 ms or more

bulk transfer: transfer multiple blocks w/o additional seek time or rotational delay, which need buffer space. Then, (that is in 1/2)

We need buffer management: e.g. LRU Strategy work for nested bp job.

When we run out of buffer space, rewrite earliest used space double buffering, 2 buffer, 1 PV read & fill other one data concurrently.

Pin-Count: # times a page required. # >

will stay in buffer.

dirty bit: page updated? 1:0.

Heap - unsorted file.

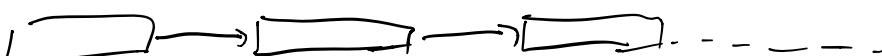
a file of data unsorted

Any record can be any block

so to find page, we need to go over all data blocks (N), on average $\frac{N}{2}$

Sorted file.

data sorted



binary search

Indices.

Primary index: specified on key field of an ordered file of records physically order the records so every records have unique value on such file.

If not unique, it is clustering index

primary or clustering but not both



key value: email Varchar(50) & pointer: 4 bytes = 54 bytes

~80% filled & 4kbs/block. \Rightarrow 3200 bytes

$3200 / 54 \approx 60$ How many pointers in 1 index block

Segment: ~60 which is each index block contains around 60 data blocks indices.

have ~200k data blocks \Rightarrow $200k / 60 \approx 3334$ index block. (sparse index) {only pick first key}

If we pick every key value on block, we will get dense

Index.

If we have $n = \#$ index blocks

Look up Time: $\log_2(n) + 1$ access the value

Dense indices take up time \rightarrow Sparse.

Secondary index: - Pointing to ~~specified on non ordering field~~ Have to be dense if unique value.
indices on ~~can be key as well~~ non-key value, which is not sync order with
key value And it has to be dense index as records are
not sorted.

So we get value of each record, pick up the key & val
in index blocks. then sort. So you can only go to a point
as data block is not sorted for this field.

Multi-level index. $\log_{\text{fanout}}(n) + 1$

index is index ...

good: fast & look up time.

bad: easy to overflow hard to adjust

B⁺-Tree.

Static Hashing (only fixed # of buckets)

Very large key space $\xrightarrow{\text{hash function}}$ addresses packed directly

good hash function:

distribute values uniformly over address space

fill buckets as much as possible

Avoid collisions

hash field space: space of elements needed to hash

B Tree: balanced, used as primary file organizations