**Carnegie Mellon**

# 15-826: Multimedia Databases and Data Mining

### Lecture #4:  Multi-key and Spatial Access Methods - I
### *C. Faloutsos*

---

**Carnegie Mellon**

# Must-Read Material

- MM-Textbook, Chapter 4
- [Bentley75] J.L. Bentley: *Multidimensional Binary Search Trees Used for Associative Searching*, CACM, 18,9, Sept. 1975.
- Ramakrinshan+Gehrke, Chapter 28.1-3

**Carnegie Mellon**

# Outline

Goal: 'Find similar / interesting things'
- Intro to DB
- Indexing - similarity search
- Data Mining

15-826                          Copyright: C. Faloutsos (2019)                          3

**Carnegie Mellon**

# Indexing - Detailed outline

- primary key indexing
- secondary key / multi-key indexing
- spatial access methods
- text
- ...

15-826                          Copyright: C. Faloutsos (2019)                          4

# Problem

- Find employees with
  - Salary in ($10K, $20K) and
  - Years-in-company in (5,7)

                                5

# Conclusions

- sec. keys: B-tree indices (+ postings lists)
- multi-key, main memory methods:
  - quad-trees
  - k-d-trees

                                6

**Carnegie Mellon**

# Sec. key indexing

- attributes w/ duplicates (eg., EMPLOYEES, with 'job-code', 'salary', 'dept')
- Query types:

15-826                          Copyright: C. Faloutsos (2019)                          7

**Carnegie Mellon**

# Sec. key indexing

- attributes w/ duplicates (eg., EMPLOYEES, with 'job-code', 'salary', 'dept')
- Query types:
  - exact match
  - partial match
    - 'job-code' = 'PGM' and 'dept' = 'R&D'
  - range queries
    - 'job-code' = 'ADMIN' and salary < 50K

15-826                          Copyright: C. Faloutsos (2019)                          8

**Carnegie Mellon**

# Sec. key indexing

- Query types - cont'd
  - boolean
    - 'job-code' = 'ADMIN' or salary>20K
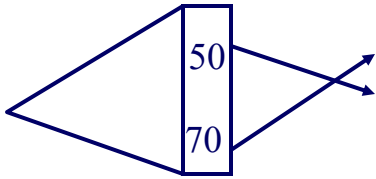  - nn
    - salary ~ 30K

**Carnegie Mellon**

# Solution?

**Carnegie Mellon**

# Solution?

- Inverted indices (usually, w/ B-trees)
- Q: how to handle duplicates?

salary-index

| Name | Job-code | Salary | Dept |
|------|----------|--------|------|
| Smith | PGM | 70 | R&D |
| Jones | ADMIN | 50 | R&D |
| …. | | | |
| Tomson | ENG | 50 | SALES |

50

70

---

**Carnegie Mellon**

# Solution

- A#1: eg., with postings lists

postings lists

salary-index

| Name | Job-code | Salary | Dept |
|------|----------|--------|------|
| Smith | PGM | 70 | R&D |
| Jones | ADMIN | 50 | R&D |
| …. | | | |
| Tomson | ENG | 50 | SALES |

50

70

**Carnegie Mellon**

# Solution

- A#2: modify B-tree code, to handle dup's

salary-index

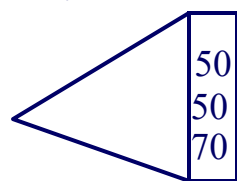| Name | Job-code | Salary | Dept |
|------|----------|--------|------|
| Smith | PGM | 70 | R&D |
| Jones | ADMIN | 50 | R&D |
| …. | | | |
| Tomson | ENG | 50 | SALES |

50
50
70

---

**Carnegie Mellon**

# How to handle Boolean Queries?

- eg., 'sal=50 AND job-code=PGM'?

salary-index

| Name | Job-code | Salary | Dept |
|------|----------|--------|------|
| Smith | PGM | 70 | R&D |
| Jones | ADMIN | 50 | R&D |
| …. | | | |
| Tomson | ENG | 50 | SALES |

50
50
70

**CarnegieMellon**

# How to handle Boolean Queries?

– from indices, find lists of qual. record-ids

– merge lists (or check real records)

salary-index

| Name | Job-code | Salary | Dept |
|------|----------|--------|------|
| Smith | PGM | 70 | R&D |
| Jones | ADMIN | 50 | R&D |
| …. | | | |
| Tomson | ENG | 50 | SALES |

50
50
70

15-826                    Copyright: C. Faloutsos (2019)                    15

---

**CarnegieMellon**

# Sec. key indexing

- easily solved in commercial DBMS:

  ```
  create index sal-index on
    EMPLOYEE (salary);
  select * from EMPLOYEE
    where salary > 50 and
    job-code = 'ADMIN'
  ```

15-826                    Copyright: C. Faloutsos (2019)                    16

**Carnegie Mellon**

# Sec. key indexing

- can create combined indices:

```
create index sj on EMPLOYEE(
  salary, job-code);
```

15-826                          Copyright: C. Faloutsos (2019)                          17

**Carnegie Mellon**

# Sec. key indexing

- can create combined indices:

```
create index sj on EMPLOYEE(
  salary, job-code);
```

Q: Drawback?

15-826                          Copyright: C. Faloutsos (2019)                          18

**Carnegie Mellon**

# Sec. key indexing

- can create combined indices:

  ```
  create index sj on EMPLOYEE(
    salary, job-code);
  ```

Q: Drawback?

A: can not answer queries on job-code

**Carnegie Mellon**

# Indexing - Detailed outline

- primary key indexing
- secondary key / multi-key indexing
    – main memory: quad-trees
    – main memory: k-d-trees
- spatial access methods
- text
- ...

# Quad-trees

- problem: find cities within 100mi from Pittsburgh
- assumption: all fit in main memory
- Q: how to answer such queries quickly?

# Quad-trees

- A: recursive decomposition of space, e.g.:

# Quad-trees

- A: recursive decomposition of space, e.g.:



PGH

PHL

ATL

10

30

SW

(30,10)

15-826          Copyright: C. Faloutsos (2019)          23

# Quad-trees

- A: recursive decomposition of space, e.g.:



PGH

PHL

ATL

20

10

30   40

SW

(30,10)

NE

(40,20)

15-826          Copyright: C. Faloutsos (2019)          24

# Quad-trees - search?

- find cities with (35<x<45, 15<y<25):

# Quad-trees - search?

- find cities with (35<x<45, 15<y<25):

**Carnegie Mellon**

# Quad-trees - search?

- pseudocode:
  range-query( tree-ptr, range)
    if (tree-ptr == NULL) exit;
    if (tree-ptr->point within range){
      print tree-ptr->point}
    for each quadrant {
      if ( range intersects quadrant ) {
        range-query( tree-ptr->quadrant-ptr, range);
      }

**Carnegie Mellon**

# Quad-trees - k-nn search?

- k-nearest neighbor algo - more complicated:
  - find 'good' neighbors and put them in a stack
  - go to the most promising quadrant, and update the stack of neighbors
  - until we hit the leaves

**Carnegie Mellon**

# Quad-trees - discussion

- great for 2- and 3-d spaces
- several variations, like fixed decomposition:

'adaptive'                                'fixed'

*z-ordering (later)*

PGH   PHL   ·   ATL

PGH   PHL   ·   ·ATL

↗ middle

---

**Carnegie Mellon**

# Quad-trees - discussion

- but: unsuitable for higher-d spaces (why?)

**Carnegie Mellon**

# Quad-trees - discussion

- but: unsuitable for higher-d spaces (why?)
- A: 2^d pointers, per node!
- Q: how to solve this problem?
- A: k-d-trees!

**Carnegie Mellon**

# Indexing - Detailed outline

- primary key indexing
- secondary key / multi-key indexing
  – main memory: quad-trees
  → – main memory: k-d-trees
- spatial access methods
- text
- ...

**Carnegie Mellon**

# k-d-trees

- Binary trees, with alternating 'discriminators'

PGH

PHL

ATL

10

30

(30,10)

SW

quad-tree

---

**Carnegie Mellon**

# k-d-trees

- Binary trees, with alternating 'discriminators'

PGH

PHL

ATL

10

30

(30,10)

W

E

**k-d-tree**

**Carnegie Mellon**

# k-d-trees

- Binary trees, with alternating 'discriminators'



ATL

(30,10)

x<=30                    x>30

PGH          PHL

ATL

10

30

---

**Carnegie Mellon**

# k-d-trees

- Binary trees, with alternating 'discriminators'



ATL

(30,10) ——— x

x<=30          x>30

PHL

(40,20) ——— y

y<=20          y>20

PGH     PHL

ATL

20

10

30 40

**CarnegieMellon**

# (Several demos/applets, e.g.)

- http://donar.umiacs.umd.edu/quadtree/points/kdtree.html

**CarnegieMellon**

# Indexing - Detailed outline

- primary key indexing
- secondary key / multi-key indexing
  - main memory: quad-trees
  - main memory: k-d-trees
    - insertion; deletion
    - range query; k-nn query
- spatial access methods
- text
- ...

# k-d-trees - insertion

- Binary trees, with alternating 'discriminators'

# k-d-trees - insertion

- discriminators: may cycle, or ....
- Q: which should we put first?

20

# k-d-trees - deletion

• How?

# k-d-trees - deletion

• Tricky! 'delete-and-promote' (or 'tombstone' = 'mark as deleted')

**Carnegie Mellon**

# k-d-trees - range query



20

10

PGH

PHL

ATL

30 40

ATL
(30,10)
x<=30          x>30
PHL
(40,20)
y<=20          y>20

**Carnegie Mellon**

# k-d-trees - range query

- similar to quad-trees: check the root; proceed to appropriate child(ren).



20

10

PGH

PHL

ATL

30 40

ATL
(30,10)
x<=30          x>30
PHL
(40,20)
y<=20          y>20

# k-d-trees - k-nn query

- e.g., 1-nn: closest city to 'X'

# k-d-trees - k-nn query

- A: check root; put in stack; proceed to child

# k-d-trees - k-nn query

- A: check root; put in stack; proceed to child



Copyright: C. Faloutsos (2019)                              47

# Indexing - Detailed outline

- primary key indexing
- secondary key / multi-key indexing
  - main memory: quad-trees
  - main memory: k-d-trees
    - insertion; deletion
    - range query; k-nn query
    - discussion
- spatial access methods
- text

Copyright: C. Faloutsos (2019)                              48

**Carnegie Mellon**

# k-d trees - discussion

- great for main memory & low 'd' (~<10)
- Q: what about high-d?
- A:
- Q: what about disk
- A:

**Carnegie Mellon**

# k-d trees - discussion

- great for main memory & low 'd' (~<10)
- Q: what about high-d?
- A: most attributes don't ever become discriminators
- Q: what about disk?
- A: Pagination problems, after ins./del. (solutions: next!)

**Carnegie Mellon**

# Conclusions

- sec. keys: B-tree indices (+ postings lists)
- multi-key, main memory methods:
  - quad-trees
  - k-d-trees

51

**Carnegie Mellon**

# References

- [Bentley75] J.L. Bentley: *Multidimensional Binary Search Trees Used for Associative Searching*, CACM, 18,9, Sept. 1975.
- [Finkel74] R.A. Finkel, J.L. Bentley: *Quadtrees: A data structure for retrieval on composite keys*, ACTA Informatica,4,1, 1974
- Applet: eg., http://donar.umiacs.umd.edu/quadtree/points/kdtree.html

52