

Vector Approximation (VA) file

- Tile d-dimensional data-space uniformly into 2^b rectangular cells.
- b bits for each approximation
- Dimension i is partitioned into 2^{b_i} partitions; this requires b_i bits:

$$b = \sum_{i=1}^d b_i.$$

Spring 2013

Searching Big data

1

Vector Approximation (VA) file

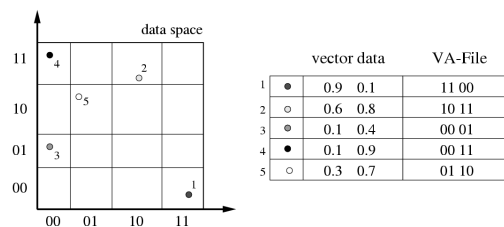
- A fixed number of bits in each dimensions (8)
- 256 partitions along each dimension
- 256^d tiles
- Approximate each point by corresponding tile
- Size of approximation = $8d$ bits = d bytes
- Size of each point = $4d$ bytes (assuming 4 bytes per dimension)

Spring 2013

Searching Big data

2

Example



Spring 2013

Searching Big data

3

Search using VA-file

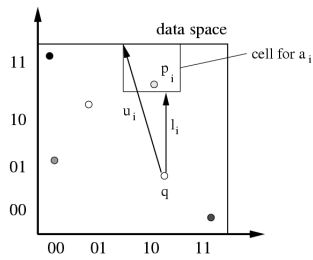
- VA file is an array of compact, geometric approximations.
- Approximations are scanned, and upper and lower bounds for points are found.
- After pruning on these bounds, remaining points are read.

Spring 2013

Searching Big data

4

Computing bounds



Spring 2013

Searching Big data

5

Simple k-NN searching

- δ = distance to kth NN so far
- For each approximation a_i
 - If $lb(q, a_i) < \delta$ then
 - Compute $r = \text{distance}(q, v_i)$
 - If $r < \delta$ then
 - Add point i to the set of NNs
 - Update δ

Spring 2013

Searching Big data

6

```

VAR ans: ARRAY OF INT; dst: ARRAY of REAL;
    a: ARRAY of Approx;  $\vec{p}$ : ARRAY of Vector;

FUNC InitCandidate(): REAL;      FUNC Candidate( $\delta$ : REAL;  $i$ : INT): REAL;
VAR j: INT;                      IF  $\delta < \text{dst}[k]$  THEN
    FOR j := 1 TO k DO             $\text{dst}[k] := \delta$ ;  $\text{ans}[k] := i$ ;
         $\text{dst}[j] := \text{MAXREAL}$ ;      SortOnDst( $\text{ans}$ ,  $\text{dst}$ ,  $k$ );
    RETURN  $\text{MAXREAL}$ ;              RETURN  $\text{dst}[k]$ ;
END-FUNC InitCandidate;          END-FUNC Candidate;

PROC VA-SSA ( $\vec{q}$ : Vector);
VAR  $i$ : INT;  $l_i$ ,  $\delta$ : REAL;
 $\delta := \text{InitCandidate}()$ ;
FOR  $i := 1$  TO  $N$  DO
     $l_i := \text{GetBounds}(a_i, \vec{v}_q)$ ;
    IF  $l_i < \delta$  THEN
         $\delta := \text{Candidate}(L_p(\vec{p}_i, \vec{q}), i)$ ;
    END-FOR;
END-PROC VA-SSA
    
```

Spring 2013

Searching Big data

Evaluation

- Advantages:
 - Simple
 - Low memory overhead
 - Everything processed sequentially: no random seeks.
- Disadvantages:
 - Performance depends on the ordering of vectors

Spring 2013

Searching Big data

8

Near-optimal NN searching

- δ = distance to kth closest $ub(q,a)$ so far
- For each approximation ai
 - Compute $lb(q,ai)$ and $ub(q,ai)$
 - If $lb(q,ai) \leq \delta$ then
 - If $ub(q,ai) < \delta$ then update δ
 - InsertHeap($lb(q,ai), i$)

/* First phase */

Spring 2013

Searching Big data

9

Near-optimal NN searching (2)

- δ = distance to kth NN so far
- Repeat
 - Examine the next entry (li, i) from the heap
 - If $\delta < li$ then break
 - else
 - Compute $r = \text{distance}(q, vi)$
 - If $r < \delta$ then
 - Add point i to the set of NNs
 - Update δ

/* Second phase */

Spring 2013

Searching Big data

10

```
PROC VA-NOA( $\vec{q}$ : Vector);
```

```
VAR  $i$ : INT;  $\delta, l_i, u_i$ : REAL; Heap: HEAP; Init(Heap);
```

```
(* PHASE - ONE *)
```

```
 $\delta := \text{InitCandidate}();$ 
```

```
FOR  $i := 1$  TO  $N$  DO
```

```
   $l_i, u_i := \text{GetBounds}(a_i, \vec{q});$ 
```

```
  IF  $l_i \leq \delta$  THEN
```

```
     $\delta := \text{Candidate}(u_i, i);$ 
```

```
    InsertHeap(Heap,  $l_i, i$ );
```

```
  END-FOR;
```

```
(* PHASE - TWO *)
```

```
 $\delta := \text{InitCandidate}();$ 
```

```
 $l_i, i := \text{PopHeap}(\text{Heap});$ 
```

```
WHILE  $l_i < \delta$  DO
```

```
   $\delta := \text{Candidate}(L_p(\vec{p}_i, \vec{q}), i);$ 
```

```
   $l_i, i := \text{PopHeap}(\text{Heap});$ 
```

```
END-WHILE;
```

```
END-PROC VA-NOA;
```

Spring 2013

Searching Big data

11

Complexity

- After first phase, 95-99% of candidates are pruned
 - Remaining number of candidates are empirically estimated to be sub-linear ($\log n$)

Spring 2013

Searching Big data

12

Experimental setup

- Data sets:
 - Synthetic: uniformly distributed data points
 - Image data: 45-dimensional 50,000 vectors
- Machine:
 - Sun SPARCstation 4
 - 85 Mhz CPU
 - 64 MB RAM

Spring 2013

Searching Big data

13

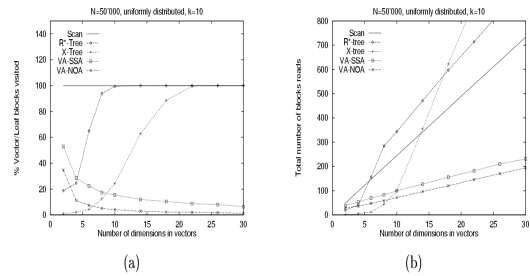


Figure 10: (a) Block selectivity and, (b) the total number of blocks accessed as a function of number of dimensions. Synthetic data set (uniformly distributed).

Spring 2013

Searching Big data

14

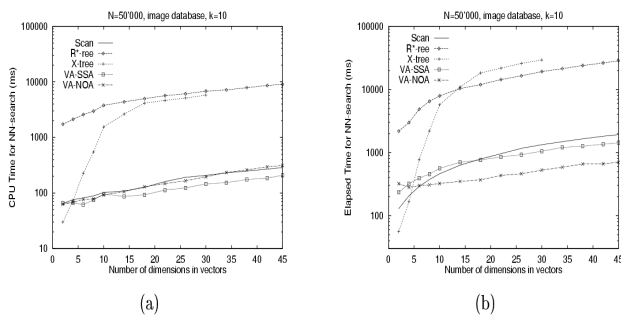


Figure 12: CPU time (a) and wall-clock time (b) for the image database

Spring 2013

Searching Big data

References

- When Is "Nearest Neighbor" Meaningful? K. Beyer, J. Goldstein, R. Ramakrishnan, U. Shaft, ICDD 1999, pp. 217-235
- A Quantitative Analysis and Performance Study for Similarity-Search Methods in High-Dimensional Spaces, R. Weber, H-J. Schek, S. Blott, VLDB 1998, pp. 194-205
- Searching in high-dimensional spaces: Index structures for improving the performance of multimedia databases, C. Böhm, S. Berchtold, D. Keim, ACM Computing Surveys, Volume 33, Issue 3, September 2001, pp. 322 - 373

Spring 2013

Searching Big data

16