


## Indexing Uncertain Categorical Data


Sarveet Singh etc.  
Presented by: Yi Zhang



## The problem

- ❖ Uncertainty in discrete-domain data
  - E.g. Nurse 10 in Room {3,5} at 10:05
  - Results from text classifiers
- ❖ Existing relational DB doesn't handle well
  - Allowing multiple values? app use own complex model to manage uncertainty
  - Just keep most likely value? lossy
  - Need proper index structure for categorical uncertain data


3/30/2008 ECE256 Course Project Yi Zhang 2



## Contributions

- ❖ Data models for uncertainty
- ❖ Semantics of queries
- ❖ Index structures
- ❖ Evaluation of index structures

3/30/2008 ECE256 Course Project Yi Zhang 3




## Data Model

- ❖ UDA: uncertain discrete attribute
  - A distribution over its discrete domain
  - Example

Employee	Department
Jim	{(Shoes, 0.5), (Sales, 0.5)}
Tom	{(Sales, 0.4), (Clothes, 0.6)}
Lin	{(Hardware, 0.6), (Sales, 0.4)}
Nancy	{(HR, 1.0)}

- UDA  $u$  can be represented as prob. vector  $\langle p_1, \dots, p_N \rangle$  s.t.  $\Pr(u=d_i) = p_i$
- For UDAs  $u$  and  $v$   $\Pr(u=v) = \sum_{i=1}^N u.p_i \times v.p_i$


3/30/2008 ECE256 Course Project Yi Zhang 4



## Distribution Similarity Metrics

- ❖ Given two UDAs, how similar are they?
- ❖ Useful for clustering in index structures
- ❖ Manhattan distance
 
$$L_1: L_1(u, v) = \sum_{i=1}^N |u.p_i - v.p_i|$$
- ❖ Euclidean distance
 
$$L_2: L_2(u, v) = \sqrt{\sum_{i=1}^N (u.p_i - v.p_i)^2}$$
- ❖ Kullback-Leibler divergence
 
$$KL(u, v): KL(u, v) = \sum_{i=1}^N u.p_i \log(u.p_i / v.p_i)$$

3/30/2008 ECE256 Course Project Yi Zhang 5



## Queries

- ❖ Equality queries
  - PEQ: given UDA  $q$ , return all tuples  $t$  from  $R$ , s.t.  $\Pr(q=t.a) \geq 0$
  - PETQ: same as PEQ except  $\Pr(q=t.a) \geq \tau$  (constant)
  - PEQ-top-k: return k tuples w/ highest equality prob.
- ❖ Distributional queries
  - DSTQ: given UDA  $q$ , divergence function  $F$ , return all tuples  $t$  from  $R$  s.t.  $F(q,t.a) \leq \tau$
  - DSQ-top-k
- ❖ Equality join queries
  - PETJ:  $R \bowtie_{R_i=S_i, \tau} S$  consists all pairs of tuples  $r, s$  from  $R, S$  s.t.  $\Pr(r.a=s.b) \geq \tau$
  - PEJ-top-k

R	a	...
S	b	...

3/30/2008 ECE256 Course Project Yi Zhang 6

### Inverted Index

- ❖ Outer list: for each domain value  $d_i$ 
  - Inner list: list of  $\langle t, p \rangle$  where  $t.d_i = p$
- ❖ Rank aggregation...
- ❖ Simple insertion/deletion
- ❖ What about searching?

3/30/2008 ECE256 Course Project Yi Zhang 7

### Search on Inverted Index

- ❖ PETQ queries: given UDA  $q$  and const  $\tau$
- ❖ Inv-index-search
  - I/O significant
- ❖ Highest-prob-first
- ❖ Row pruning
- ❖ Column pruning
  - Needs random I/O

3/30/2008 ECE256 Course Project Yi Zhang 8

### PDR-tree

- ❖ Consider UDA as point in high-dim space
- ❖ Similar UDAs stored in a page
- ❖ Page described by MBR
- ❖ Parent's MBR covers all its children

3/30/2008 ECE256 Course Project Yi Zhang 10

### Operations on PDR-trees

- ❖ Insert( $u$ )
  - To insert into current page, update MBR; then choose best child page and insert.
    - Min area increase
    - Most similar MBR
  - Start from root and recurse
- ❖ Split()
  - Top-down: pick two children MBRs as cluster centers and all other UDAs go to a closer cluster
  - Bottom-up: each UDA forms an independent cluster and proceed by merging closest pair of clusters until two remains

3/30/2008 ECE256 Course Project Yi Zhang 10

### Queries on PDR-trees

- ❖ PETQ( $q, \tau$ )
  - Depth-first search
  - Enters an node  $c$  if qualifies  $c.mbr \cdot q \geq \tau$ ; prune otherwise
  - At leaf level, check each UDA and output
  - For top-k queries, update threshold dynamically during search
  - Doesn't work with distribution similarity queries...

3/30/2008 ECE256 Course Project Yi Zhang 11

### Experimental Setup

- ❖ Datasets
  - Real
    - CRM databases, 100k text doc, 50 categories
      - CRM1: prob values generated by auto categorization
      - CRM2: unsupervised fuzzy clustering
  - Synthetic
    - Uniform: 5 categories, prob. random chosen, 10k tuples
    - Pairwise: 5 categories, each tuple has 2 non-zero categories w/ roughly equal prob., 10k tuples
    - Gen3: for each tuple, size of non-zero categories follows geometric dist.
- ❖ Page size 8KB, 100 blocks for buffering
- ❖ PETQ and PEQ-top-k queries
- ❖ Measure # of I/O operations

3/30/2008 ECE256 Course Project Yi Zhang 12

### Divergence Measures

- ❖ KL: if hit an MBR, most UDAs will qualify cause they are similar in dist.
- ❖ Top-k: explore more tuples

**Figure 4. L1 vs L2 vs KL (PDR-tree)**

3/30/2008 ECE256 Course Project Yi Zhang 13

### Inverted Index vs. PDR-tree

- ❖ Non-zero probability in many categories
- ❖ Large number of lists accessed in inverted index

**Figure 6. Inverted Index vs PDR-tree (CRM1)**

3/30/2008 ECE256 Course Project Yi Zhang 14

### Data/Domain Size

- ❖ Inverted Index
  - Reduction in avg length of each list as # of lists increase
- ❖ PDR-tree
  - Data generation related
  - Non-zero entries at both ends of experimental space smaller than in the middle

**Figure 8. Scalability with Size of Data**

**Figure 9. Scalability with Domain Size**

3/30/2008 ECE256 Course Project Yi Zhang 15

### Discussion

- ❖ How to use the mentioned criteria to insert a UDA into a PDR-tree
- ❖ No motivation for MBR-based approach used in PDR-Tree
- ❖ Not obvious how to extend PDR-Tree to work for Distributional Similarity Queries

3/30/2008 ECE256 Course Project Yi Zhang 16